

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

APPLICATION FOR LETTERS PATENT

**EV 355228348**

**Network Load Balancing  
with Traffic Routing**

Inventor(s):

**Christopher L. Darling**

**Sean B. House**

**Aamer Hydrie**

**Joseph M. Joy**

**Robert V. Welland**

ATTORNEY'S DOCKET NO. **MS1-1518US**

# Network Load Balancing with Traffic Routing

## RELATED PATENT APPLICATIONS

This U.S. Nonprovisional Application for Letters Patent (i) is a continuation-in-part of co-pending U.S. Nonprovisional Application for Letters Patent No. 10/610,506 (filed June 30, 2003), (ii) is a continuation-in-part of co-pending U.S. Nonprovisional Application for Letters Patent No. 10/610,519 (filed June 30, 2003), and (iii) is a continuation-in-part of co-pending U.S. Nonprovisional Application for Letters Patent No. 10/610,321 (filed June 30, 2003).

Specifically, this U.S. Nonprovisional Application for Letters Patent is a continuation-in-part of, and hereby incorporates by reference herein the entire disclosure of, co-pending U.S. Nonprovisional Application for Letters Patent No. 10/610,506, filed June 30, 2003, and entitled "Flexible Network Load Balancing".

Specifically, this U.S. Nonprovisional Application for Letters Patent is also a continuation-in-part of, and hereby incorporates by reference herein the entire disclosure of, co-pending U.S. Nonprovisional Application for Letters Patent No. 10/610,519, filed June 30, 2003, and entitled "Network Load Balancing with Host Status Information".

Specifically, this U.S. Nonprovisional Application for Letters Patent is also a continuation-in-part of, and hereby incorporates by reference herein the entire disclosure of, co-pending U.S. Nonprovisional Application for Letters Patent No.

1 10/610,321, filed June 30, 2003, and entitled "Network Load Balancing with  
2 Session Information".

3  
4 **TECHNICAL FIELD**

5 This disclosure relates in general to network load balancing and in  
6 particular, by way of example but not limitation, to network load balancing with  
7 traffic routing and the high availability thereof.

8  
9 **BACKGROUND**

10 Communication, and many facets of life that involve communication, has  
11 been greatly impacted by the Internet. The Internet enables information to be  
12 communicated between two people and/or entities quickly and relatively easily.  
13 The Internet includes many network nodes that are linked together such that  
14 information may be transferred between and among them. Some network nodes  
15 may be routers that propagate a packet from one link to another, may be individual  
16 client computers, may be personal networks for different entities (e.g., intranets  
17 for businesses), and so forth.

18 For this personal network case, as well as others, packets arriving at an  
19 Internet node or nodes are distributed to other nodes within the personal network.  
20 Such a personal network may be formed, for example, from a set of servers that  
21 can each work on packets that arrive at the personal network. A business, a  
22 university, a government office, etc. may receive many packets in a short  
23 timeframe at its personal network. In order to respond in a timely manner and to  
24 reduce the likelihood of rejection or loss of arriving packets, the personal network  
25

1 may rely on multiple servers that can each work on the arriving packets  
2 simultaneously.

3 The arriving packets are often inquiries pertaining to certain information,  
4 such as a document, a catalog item, a web page, and so forth. The arriving packets  
5 can also pertain to an economic transaction between a customer and a merchant.  
6 Other purposes for the packets of a packet-based communication are possible.  
7 Regardless, the arriving packets are distributed among different servers of a set of  
8 servers to accommodate a rapid arrival of the packets and/or complex  
9 communication exchanges.

10 The distribution of arriving packets among different servers of a set of  
11 servers is often termed network load balancing. In other words, a load balancing  
12 operation may be performed on packets as they arrive at a node or nodes of the  
13 Internet when the node or nodes constitute a personal network and/or when they  
14 connect the personal network to the Internet.

15 Such a load balancing operation is accomplished using dedicated hardware  
16 that fronts the personal network at the node or nodes that connect the personal  
17 network to the Internet and/or that provide a presence for the personal network on  
18 the Internet. The physical hardware that performs the load balancing operation is  
19 usually duplicated in its entirety to realize redundancy and improve availability of  
20 the load balancing operation. To increase capacity for load balancing operations,  
21 more-powerful hardware that replicates the entirety of the previous load balancing  
22 hardware, and thus the operational capability thereof, is substituted for the  
23 previous load balancing hardware. Such scaling up of the load balancing  
24 operational capabilities is therefore confined to increasing the power of the  
25 hardware via substitution thereof.



1 To implement a load balancing operation, the hardware usually performs a  
2 round robin distribution of arriving connection requests. In other words, arriving  
3 connection requests are distributed to servers of a set of servers in a linear,  
4 repeating manner with a single connection request being distributed to each server.  
5 This round-robin load balancing distribution of connections is typically utilized  
6 irrespective of the condition of the personal network or the nature of an arriving  
7 connection request. If a load balancing operation does extend beyond a round  
8 robin distribution, these other factors are only considered to the extent that they  
9 may be inferred from network traffic and/or from a congestion level of the  
10 personal network.

11 Accordingly, there is a need for schemes and/or techniques that improve  
12 network load balancing and/or the options associated therewith.

### 13 14 **SUMMARY**

15 In an exemplary method implementation, a method includes: receiving a  
16 packet requesting a new connection at a forwarding component; sending the  
17 packet from the forwarding component to a classifying component; selecting, by  
18 the classifying component, a route for the new connection; and plumbing, by the  
19 classifying component, the route for the new connection by causing a new entry to  
20 be added in a local routing table of the forwarding component.

21 In an exemplary media implementation, one or more processor-accessible  
22 media include processor-executable instructions that, when executed, enable a  
23 system to perform actions including: receiving a first packet for a connection at  
24 first forwarding functionality; plumbing a route for the connection at the first  
25 forwarding functionality; receiving a second packet for the connection at second

1 forwarding functionality; and plumbing the route for the connection at the second  
2 forwarding functionality using a distributed session tracking table.

3 In an exemplary system implementation, a system includes: a forwarding  
4 component that forwards packets; a classifying component that classifies packets  
5 and is capable of classifying packets for the forwarding component; a session  
6 tracking component that tracks sessions for at least one of the forwarding  
7 component and the classifying component; a health and load handling component  
8 that is capable of providing health and load information to the classifying  
9 component; and a high availability mechanism that provides detection of, handling  
10 of, and recovery from a failure of one or more of the forwarding component, the  
11 classifying component, the session tracking component, and the health and load  
12 handling component.

13 In another exemplary media implementation, one or more processor-  
14 accessible media include processor-executable instructions that, when executed,  
15 direct a system to perform actions including: receiving a token allotment at traffic  
16 routing functionality from health and load functionality, the token allotment  
17 having a first multiplicity of tokens for a first destination and a second multiplicity  
18 of tokens for a second destination; consuming, by the traffic routing functionality,  
19 a token of the first multiplicity of tokens when selecting the first destination for a  
20 connection request; and consuming, by the traffic routing functionality, a token of  
21 the second multiplicity of tokens when selecting the second destination for a  
22 connection request.

23 In yet another exemplary media implementation, one or more processor-  
24 accessible media include processor-executable instructions for load balancing  
25 infrastructure that, when executed, enable a system to perform actions including:

1 establishing a first connection with a client; receiving a first request from the  
2 client via the first connection; determining, responsive to session information  
3 and/or health and load information, that the first request is to be routed to a first  
4 host via a second connection; receiving a second request from the client via the  
5 first connection; and determining, responsive to the session information and/or the  
6 health and load information, that the second request is to be routed to a second  
7 host via a third connection.

8 Other method, system, approach, apparatus, application programming  
9 interface (API), device, media, procedure, arrangement, etc. implementations are  
10 described herein.

## 11 12 **BRIEF DESCRIPTION OF THE DRAWINGS**

13 The same numbers are used throughout the drawings to reference like  
14 and/or corresponding aspects, features, and components.

15 FIG. 1 is an exemplary network load balancing paradigm that illustrates a  
16 load balancing infrastructure and multiple hosts.

17 FIG. 2 is an exemplary network load balancing paradigm that illustrates  
18 multiple load balancing units and multiple hosts.

19 FIG. 3 illustrates an exemplary load balancing unit having separated  
20 functionality and an exemplary host.

21 FIG. 4 illustrates exemplary network load balancing infrastructure having  
22 separated classifying and forwarding functionality.

23 FIG. 5 is a flow diagram that illustrates an exemplary method for scaling  
24 out network load balancing infrastructure into different configurations.  
25

1        FIG. 6 illustrates a first exemplary network load balancing infrastructure  
2 configuration from a device perspective.

3        FIG. 7 illustrates a second exemplary network load balancing infrastructure  
4 configuration from a device perspective.

5        FIGS. 8A and 8B illustrate first and second exemplary network load  
6 balancing infrastructure configurations from a component perspective.

7        FIGS. 9A and 9B illustrate first and second exemplary network load  
8 balancing infrastructure configurations from a resource perspective.

9        FIG. 10 illustrates an exemplary network load balancing approach that  
10 involves host status information.

11       FIG. 11 is a flow diagram that illustrates an exemplary method for network  
12 load balancing that involves host status information.

13       FIG. 12 illustrates an exemplary network load balancing approach that  
14 involves health and load information.

15       FIG. 13A is an exemplary health and load table as illustrated in FIG. 12.

16       FIG. 13B is an exemplary consolidated health and load cache as illustrated  
17 in FIG. 12.

18       FIG. 14 is a flow diagram that illustrates an exemplary method for network  
19 load balancing that involves health and load information.

20       FIG. 15 illustrates an exemplary message protocol for communications  
21 between the hosts and load balancing units that are illustrated in FIG. 12.

22       FIG. 16 illustrates an exemplary message transmission scheme for  
23 communications between the hosts and load balancing units that are illustrated in  
24 FIG. 12.

1           FIGS. 17A and 17B illustrate exemplary health and load information proxy  
2 storage scenarios for health and load tables of FIG. 13A and for consolidated  
3 health and load caches of FIG. 13B, respectively.

4           FIG. 18 illustrates an exemplary target host allotment procedure that utilizes  
5 health and load information.

6           FIG. 19 illustrates an exemplary network load balancing approach that  
7 involves session information.

8           FIG. 20 illustrates an exemplary network load balancing approach that  
9 involves communicating session information using notifications and messages.

10          FIG. 21 is a flow diagram that illustrates an exemplary method for network  
11 load balancing that involves communicating session information using  
12 notifications and messages.

13          FIG. 22 illustrates an exemplary approach to managing session information  
14 at multiple load balancing units.

15          FIG. 23A is an exemplary session table as illustrated in FIG. 20.

16          FIG. 23B is an exemplary distributed atom manager (DAM) table (DAMT)  
17 as illustrated in FIG. 22.

18          FIG. 24 is a flow diagram that illustrates an exemplary method for  
19 managing session information at multiple load balancing units.

20          FIG. 25 illustrates exemplary network load balancing infrastructure having  
21 request routing functionality.

22          FIG. 26 is a flow diagram that illustrates an exemplary method for routing  
23 incoming packets with regard to (i) session information and (ii) health and load  
24 information.

1        FIG. 27 illustrates an exemplary traffic routing flow in the absence of  
2 failures.

3        FIG. 28 illustrates an exemplary traffic routing flow in the presence of  
4 failure(s).

5        FIG. 29 illustrates additional exemplary failover procedures for high  
6 availability of network load balancing infrastructure.

7        FIG. 30 illustrates an exemplary operational implementation of traffic  
8 routing interaction with health and load information.

9        FIG. 31 illustrates exemplary high availability mechanisms for network  
10 load balancing infrastructure.

11       FIG. 32 illustrates an exemplary approach to application-level network load  
12 balancing with connection migration.

13       FIG. 33 is a flow diagram that illustrates an exemplary method for  
14 migrating a connection from a first device to a second device.

15       FIG. 34 illustrates an exemplary approach to connection migration from the  
16 perspective of an originating device.

17       FIG. 35 illustrates an exemplary approach to connection migration from the  
18 perspective of a targeted device.

19       FIG. 36 illustrates an exemplary approach to an offloading procedure for a  
20 connection migration.

21       FIG. 37 illustrates an exemplary approach to an uploading procedure for a  
22 connection migration.

23       FIG. 38 illustrates an exemplary approach to packet tunneling between a  
24 forwarder and a host.

1 FIG. 39 is a flow diagram that illustrates an exemplary method for packet  
2 tunneling between a first device and a second device.

3 FIG. 40 illustrates an exemplary computing (or general device) operating  
4 environment that is capable of (wholly or partially) implementing at least one  
5 aspect of network load balancing as described herein.

## 6 7 **DETAILED DESCRIPTION**

### 8 **Exemplary Network Load Balancing Paradigms**

9 This section describes exemplary paradigms for network load balancing  
10 and is used to provide foundations, environments, contexts, etc. for the  
11 descriptions in the following sections. This section primarily references FIGS. 1-  
12 3.

13 FIG. 1 is an exemplary network load balancing paradigm 100 that illustrates  
14 a load balancing infrastructure 106 and multiple hosts 108. Exemplary network  
15 load balancing paradigm 100 includes multiple clients 102(1), 102(2) ... 102(m)  
16 and multiple hosts 108(1), 108(2) ... 108(n), as well as network 104 and load  
17 balancing infrastructure 106.

18 Each of clients 102 may be any device that is capable of network  
19 communication, such as a computer, a mobile station, an entertainment appliance,  
20 another network, and so forth. Clients 102 may also relate to a person and/or  
21 entity that is operating a client device. In other words, clients 102 may comprise  
22 logical clients that are users and/or machines. Network 104 may be formed from  
23 one or more networks, such as the Internet, an intranet, a wired or wireless  
24 telephone network, and so forth. Additional examples of devices for clients 102  
25 and network types/topologies for network 104 are described below with reference

1 to FIG. 40 in the section entitled “Exemplary Operating Environment for  
2 Computer or Other Device”.

3 Individual clients 102 are capable of communicating with one or more  
4 hosts 108, and vice versa, across network 104 via load balancing infrastructure  
5 106. Hosts 108 host one or more applications for interaction/communication with  
6 clients 102, for use by clients 102, and so forth. Each host 108 may correspond to  
7 a server and/or a device, multiple servers and/or multiple devices, part of a server  
8 and/or part of a device, some combination thereof, and so forth. Particular  
9 implementations for hosts 108 are described further below in the context of  
10 different network load balancing situations. (However, back-end support for hosts  
11 108 is generally not shown for the sake of clarity.) Furthermore, additional  
12 examples of devices for hosts 108 are also described below with reference to FIG.  
13 40 in the section entitled “Exemplary Operating Environment for Computer or  
14 Other Device”.

15 Load balancing infrastructure 106 is reachable or locatable through network  
16 104 at one or more virtual internet protocol (IP) addresses. Communications from  
17 clients 102 (or other nodes) that are directed to the virtual IP address of load  
18 balancing infrastructure 106 are received there and forwarded to a host 108. Load  
19 balancing infrastructure 106 is comprised of hardware and/or software  
20 components (not explicitly shown in FIG. 1).

21 Although load balancing infrastructure 106 is shown as an integral ellipse,  
22 the infrastructure to effectuate load balancing may also be distributed to other  
23 aspects of exemplary network load balancing paradigm 100. For example,  
24 software component(s) of load balancing infrastructure 106 may be located at one  
25 or more of hosts 108 as is described further below. Examples of architectures for



1 load balancing infrastructure 106 are described below with reference to FIG. 40 in  
2 the section entitled "Exemplary Operating Environment for Computer or Other  
3 Device".

4 As indicated at (1), one or more of hosts 108 may provide host status  
5 information from hosts 108 to load balancing infrastructure 106. This host status  
6 information may be application specific. Examples of such host status information  
7 are described further below and include health and/or load information, session  
8 information, etc. for hosts 108. A particular implementation that includes  
9 providing health and/or load information from hosts 108 to load balancing  
10 infrastructure 106 is described below in the section entitled "Exemplary Health  
11 and Load Handling".

12 At (2), a request is sent from client 102(1) across network 104 to load  
13 balancing infrastructure 106 at the virtual IP address thereof. The content, format,  
14 etc. of a request from a client 102 may depend on the application to which the  
15 request is directed, and the term "request" may implicitly include a response or  
16 responses from host(s) 108, depending on the context. Kinds of client requests  
17 include, but are not limited to:

- 18 1. Hyper text transfer protocol (HTTP) GET requests  
19 from a client using a browser program. Depending on the  
20 application (and more specifically, on the uniform resource locator  
21 (URL) of the requests), it may be better to service the requests by  
22 different sets of hosts, and the existence of a client "session" state on  
23 the hosts may militate that requests from specific clients be routed to  
24 specific hosts. The requests may be over a secure sockets layer  
25 (SSL) (or other encrypted) connection.

1           2.     Virtual private network (VPN) connections (e.g., the  
2     hosts are a set of VPN servers). In this case, the “request” can be  
3     considered to be a layer-2 tunneling protocol (L2TP) or point-to-  
4     point tunneling protocol (PPTP) “connection” (the latter is a  
5     combination of a transmission control protocol (TCP) control  
6     connection and associated generic routing encapsulation (GRE) data  
7     traffic).

8           3.     Terminal server connections (e.g., the hosts are a set of  
9     terminal servers).

10          4.     Proprietary requests in the form of individual TCP  
11     connections (one per request) employing a proprietary application-  
12     specific protocol.

13          5.     Simple object access protocol (SOAP) requests.

14          6.     Real-time communication requests involving control  
15     information over a TCP connection and latency-sensitive media  
16     streaming over real-time protocol (RTP).

17     Thus, requests can take many diverse, application-specific forms. In certain  
18     described implementations, load balancing infrastructure 106 may make  
19     application-specific forwarding decisions.

20           At (3), load balancing infrastructure 106 forwards the request from 102(1)  
21     to host 108(2) (in this example). Load balancing infrastructure 106 may consider  
22     one or more of many factors when selecting a host 108 to which the request is to  
23     be forwarded, depending on which implementation(s) described herein are being  
24     employed. For example, load balancing infrastructure 106 may take into account:  
25     the application health and/or load information of each host 108, session  
   information relating to client 102(1) as stored at a host 108, and so forth.

1 FIG. 2 is an exemplary network load balancing paradigm 200 that illustrates  
2 multiple load balancing units 106 and multiple hosts 108. Specifically, load  
3 balancing infrastructure 106 is shown as multiple load balancing units 106(1),  
4 106(2) ... 106(u) in exemplary network load balancing paradigm 200.  
5 Additionally, two router and/or switches 202(1) and 202(2) are illustrated.

6 Router/switches 202, if present, may be considered as part of or separate  
7 from load balancing infrastructure 106 (of FIG. 1). Router/switches 202 are  
8 responsible for directing overall requests and individual packets that are received  
9 from network 104 to the shared virtual IP (VIP) address(es) of load balancing units  
10 106. If a first router/switch 202 fails, the second router/switch 202 may takeover  
11 for the first. Although two router/switches 202 are illustrated, one or more than  
12 two router/switches 202 may alternatively be employed.

13 Router/switches 202 may be ignorant of the load balancing infrastructure or  
14 load-balancing aware. If router/switches 202 are not load-balancing aware, one of  
15 two exemplary options may be employed: For a first option, one load balancing  
16 unit 106 is "assigned" the shared VIP address, and all network traffic is forwarded  
17 thereto. This one load balancing unit 106 then evenly redistributes the traffic  
18 across the other load balancing units 106. However, there are bottleneck and  
19 failover issues with this first option (which can be mitigated if multiple VIP  
20 addresses are shared and are split between multiple load balancing units 106). For  
21 a second option, router/switches 202 are "tricked" into directing network traffic to  
22 all load balancing units 106, which individually decide what traffic each should  
23 accept for load balancing. However, there are inefficient effort duplication and  
24 switch performance/compatibility issues with this second option.  
25

1        If, on the other hand, router/switches 202 are load-balancing aware,  
2 router/switches 202 can be made to distribute incoming network traffic  
3 between/among multiple load balancing units 106 (e.g., in a round-robin fashion).  
4 It should be understood that such load-balancing-aware routers/switches 202 are  
5 capable of performing load balancing functions at a rudimentary level (e.g., in  
6 hardware). For example, load-balancing-aware routers/switches 202 can perform  
7 simple IP-address-based session affinity so that all packets from a specific source  
8 IP address are directed to a same load balancing unit 106.

9        Each separately-illustrated load balancing unit 106 of load balancing units  
10 106 may represent one physical device, multiple physical devices, or part of a  
11 single physical device. For example, load balancing unit 106(1) may correspond  
12 to one server, two servers, or more. Alternatively, load balancing unit 106(1) and  
13 load balancing unit 106(2) may together correspond to a single server. An  
14 exemplary load balancing unit 106 is described further below from a functional  
15 perspective with reference to FIG. 3.

16        Two exemplary request paths [1] and [2] are illustrated in FIG. 2. For  
17 request path [1], client 102(2) transmits a request over network 104 that reaches  
18 router/switch 202(1). Router/switch 202(1) directs the packet(s) of the request that  
19 originated from client 102(2) to load balancing unit 106(1). Load balancing unit  
20 106(1) then forwards the packet(s) of the request to host 108(1) in accordance with  
21 some load-balancing functionality (e.g., policy). For request path [2], client  
22 102(m) transmits a request over network 104 that reaches router/switch 202(2).  
23 Router/switch 202(2) directs the packet(s) of the request that originated from  
24 client 102(m) to load balancing unit 106(u). Load balancing unit 106(u) then  
25 forwards the packet(s) of the request to host 108(n) in accordance with some load-

1 balancing functionality. Exemplary load-balancing functionality is described  
2 further below with reference to FIG. 3.

3 FIG. 3 illustrates an exemplary load balancing unit 106 having separated  
4 functionality and an exemplary host 108. Load balancing unit 106 includes seven  
5 (7) functional blocks 302-314. These functional blocks of load balancing unit 106  
6 may be realized at least partially using software. Host 108 includes one or more  
7 applications 316. In a described implementation, load balancing unit 106 includes  
8 a forwarder 302, a classifier 304, a request router 306, a session tracker 308, a  
9 connection migrator 310, a tunneler 312, and a health and load handler 314.

10 Health and load handler 314 is located partly at hosts 108 and partly on  
11 devices of load balancing units 106. Health and load handler 314 monitors the  
12 health and/or load (or more generally the status) of hosts 108 so that health and/or  
13 load information thereof may be used for the load-balancing functionality (e.g.,  
14 when making load-balancing decisions). Exemplary implementations for health  
15 and load handler 314 are described further below, particularly in the section  
16 entitled "Exemplary Health and Load Handling".

17 Session tracker 308 may also be located partly at hosts 108 and partly on  
18 devices of load balancing units 106. Session tracker 308 monitors sessions that  
19 are established by clients 102 so that reconnections/continuations of previously-  
20 established sessions may be facilitated by the load-balancing functionality. For  
21 example, some applications keep application-specific client session data on the  
22 hosts (which is also a type of host status information). These applications  
23 typically expect that clients use the same host for the duration of any given  
24 session. Exemplary types of sessions include: (i) a TCP connection (which is,  
25

1 strictly speaking, a session); (ii) an SSL session; (iii) a secure IP (IPsec) session;  
2 (iv) an HTTP cookie-based session; and so forth.

3 Although session tracker 308 is illustrated as a discrete block in load  
4 balancing unit 106, session tracking functionality of session tracker 308 may  
5 actually be implemented at a global level. In other words, session affinity is  
6 supported across multiple load balancing units 106. Session tracker 308 includes a  
7 centralized database and/or a distributed database of session information in order  
8 to preserve session affinity. Exemplary implementations for session tracker 308,  
9 with an emphasis on a distributed database approach, are described further below,  
10 particularly in the section entitled “Exemplary Session Tracking”.

11 Classifier 304 uses the data acquired and maintained by health and load  
12 handler 314 and/or session tracker 308, possibly in conjunction with other factors,  
13 to classify incoming requests. In other words, classifier 304 selects a target host  
14 108 for each incoming request from a client 102. Forwarder 302 forwards client  
15 requests (and/or the packets thereof) in accordance with the targeted host 108 as  
16 selected by classifier 304. Forwarder 302 and classifier 304 may operate on a per-  
17 packet basis. Exemplary implementations for forwarder 302 and classifier 304 are  
18 described further below, particularly in the sections entitled “Exemplary Approach  
19 to Flexible Network Load Balancing” and “Exemplary Classifying, Forwarding,  
20 and Request Routing”.

21 Request router 306, as contrasted with per-packet implementations of  
22 forwarder 302 and classifier 304, can act as a proxy for an application running on  
23 a host 108. For example, request router 306 may terminate TCP connections,  
24 parse (perhaps partially) each logical request from a client 102, and resubmit each  
25 logical request to the targeted host 108. Consequently, each logical request from a

1 client 102 may be directed to a different host 108, depending on the decisions  
2 made by request router 306. Furthermore, request router 306 may perform pre-  
3 processing on a connection (e.g., SSL decryption), may choose to absorb certain  
4 requests (e.g., because request router 306 maintains a cache of responses), may  
5 arbitrarily modify requests before forwarding them to hosts 108, and so forth.  
6 Exemplary implementations for request router 306 are also described further  
7 below, particularly in the sections entitled “Exemplary Approach to Flexible  
8 Network Load Balancing” and “Exemplary Classifying, Forwarding, and Request  
9 Routing”.

10 Connection migrator 310 enables a connection to be initially terminated at  
11 load balancing unit 106 and then migrated such that the connection is subsequently  
12 terminated at host 108. This connection migration can facilitate application-level  
13 load balancing. Connection migrator 310 is capable of migrating a connection  
14 from load balancing unit 106 to a host 108 in such a manner that that the original  
15 termination at load balancing unit 106 is transparent to a requesting client 102 and  
16 to applications 316 of the newly-terminating host 108. Tunneler 312 may utilize  
17 an encapsulation scheme for the tunneling of packets that does not introduce an  
18 overhead to each tunneled packet.

19 The functionality of tunneler 312 may also be used in situations that do not  
20 involve a connection migration. Furthermore, connection migrator 310 and/or  
21 tunneler 312 may additionally be used in non-load-balancing implementations.  
22 Exemplary implementations for connection migrator 310, as well as for tunneler  
23 312, are described further below, particularly in the section entitled “Exemplary  
24 Connection Migrating with Optional Tunneling and/or Application-Level Load  
25 Balancing”.

1 Any given implementation of a load balancing unit 106 may include one or  
2 more of the illustrated functions. Although illustrated separately, each of the  
3 functions of blocks 302-314 may actually be interrelated with, overlapping with,  
4 and/or inclusive of other functions. For example, health and/or load information  
5 of health and load handler 314 may be used by classifier 304. Also, connection  
6 migrator 310 and tunneler 312 work in conjunction with forwarder 302 and  
7 classifier 304. Certain other exemplary overlapping and interactions are described  
8 herein below.

9 In a described implementation, host 108 runs and provides access to one or  
10 more applications 316. Generally, applications 316 include file delivery programs,  
11 web site management/server programs, remote access programs, electronic mail  
12 programs, database access programs, and so forth. Specifically, applications 316  
13 may include, but are not limited to, web servers such as Internet Information  
14 Server<sup>®</sup> (IIS) from Microsoft<sup>®</sup> Corporation, terminal servers such as Microsoft<sup>®</sup>  
15 Terminal Server<sup>™</sup>, and firewall and proxy products such as Internet Security and  
16 Acceleration Server<sup>™</sup> (ISA). Although the specific application 316 examples in  
17 the preceding sentence relate to Microsoft<sup>®</sup> products, network load balancing as  
18 described herein is not limited to any particular vendor(s), application(s), or  
19 operating system(s).

#### 20 **Exemplary Approach to Flexible Network Load Balancing**

21 This section illuminates how the network load balancing implementations  
22 described in this and other sections herein provide a flexible approach to network  
23 load balancing. This section primarily references FIGS. 4-9B.

24 As noted above, network load balancing functionality may be scaled up by  
25 replacing a first network load balancer with a second, bigger and more powerful



1 network load balancer. The hardware capabilities of the second network load  
2 balancer replicate the entirety of the hardware capabilities of the first network load  
3 balancer, except that a greater capacity is provided. This is an inflexible approach  
4 that can be very inefficient, especially when only one network load balancing  
5 feature is limiting performance and precipitating an upgrade of a network load  
6 balancer.

7 FIG. 4 illustrates exemplary network load balancing infrastructure having  
8 separated classifying and forwarding functionality. The separated classifying  
9 functionality and forwarding functionality are represented by classifier 304 and  
10 forwarder 302, respectively. Although classifying and forwarding functions are  
11 described further below, especially in the section entitled "Exemplary Classifying,  
12 Forwarding, and Request Routing", an initial description is presented here as an  
13 example of interaction between network load balancing infrastructure functionality  
14 and hosts 108.

15 In a described implementation, forwarder 302 corresponds to, and is the  
16 network endpoint for, the virtual IP (VIP) address (or addresses). Forwarder 302  
17 is a relatively low-level component that makes simplified and/or elementary policy  
18 decisions, if any, when routing packets to a further or final destination. Forwarder  
19 302 consults a routing table to determine this destination. Classifier 304 populates  
20 the routing table based on one or more factors (e.g., host status information),  
21 which are described further in other sections herein.

22 Clients 102 and hosts 108 also correspond to indicated network addresses.  
23 Specifically, client 102(1) corresponds to address C1, client 102(2) corresponds to  
24 address C2 ... client 102(m) corresponds to address Cm. Also, host 108(1)  
25

1 corresponds to address H1, host 108(2) corresponds to address H2 ... host 108(n)  
2 corresponds to address Hn.

3 Five communication paths (1)-(5) are shown in FIG. 4. Communication  
4 path (1) is between client 102(1) and forwarder 302, and communication path (5)  
5 is between forwarder 302 and host 108(1). Communication paths (2)-(4) are  
6 between forwarder 302 and classifier 304. For simplicity in this example, the  
7 connection associated with communication paths (1)-(5) is an HTTP TCP  
8 connection. Furthermore, load balancing in this example relates to routing  
9 incoming connections to the least loaded host 108, at least without any explicit  
10 consideration of application-level load balancing.

11 Communication paths (1)-(5) indicate how forwarder 302 and classifier 304  
12 load-balance a single HTTP TCP connection from client 102(1). At (1), client  
13 102(1) initiates the TCP connection by sending a TCP SYN packet addressed to  
14 the VIP address. The routing infrastructure of network 104 routes this packet to  
15 forwarder 302 via router/switch 202(1), which is the "closest" router/switch 202 to  
16 forwarder 302.

17 At (2), forwarder 302 consults a routing table, which may be internal to  
18 forwarder 302 or otherwise accessible therefrom, in order to look up this  
19 connection. This connection may be identified in the routing table by the TCP/IP  
20 4-tuple (i.e., source IP address, source TCP port, destination IP address, destination  
21 TCP port). Because this is the first packet of the connection, there is no entry in  
22 the routing table. Forwarder 302 therefore applies the "default route" action,  
23 which is to send this packet to classifier 304.

24 At (3), classifier 304 consults its (e.g., consolidated) cache of host status  
25 information for hosts 108(1), 108(2) ... 108(n). Classifier 304 concludes that host

1 108(1) is available and the least loaded host 108 at this instant for this example.  
2 Classifier 304 also “plumbs” a route in the routing table consulted by forwarder  
3 302 for this TCP connection. For example, classifier 304 adds a route entry or  
4 instructs forwarder 302 to add a route entry to the routing table that maps the TCP  
5 connection (e.g., identified by the TCP 4-tuple) to a specific destination host 108,  
6 which is host 108(1) in this example. More particularly, the route entry specifies  
7 the network address H1 of host 108(1).

8 At (4), classifier 304 sends the TCP SYN packet back to forwarder 302.  
9 Alternatively, classifier 304 may forward this initial TCP SYN packet to host  
10 108(1) without using forwarder 302. Other options available to classifier 304 are  
11 described further below.

12 At (5), forwarder 302 can access a route entry for the connection  
13 represented by the SYN packet, so it forwards the packet to host 108(1) at address  
14 H1. Forwarder 302 also forwards all subsequent packets from client 102(1) for  
15 this connection directly to host 108(1). In other words, forwarder 302 can avoid  
16 further interaction with classifier 304 for this connection. One or a combination of  
17 mechanisms, which are described further below, may be used to delete the route  
18 entry when the connection ceases.

19 For communication path (5) in many protocol environments, forwarder 302  
20 cannot simply send the packets from client 102(1) as-is to host 108(1) at network  
21 address H1 because these packets are addressed to the VIP address, which is  
22 hosted by forwarder 302 itself. Instead, forwarder 302 may employ one or more  
23 of the following exemplary options:

- 24 1. Forwarder 302 performs Network Address Translation  
25 (NAT) by (i) overwriting the source (client 102(1)) IP address (C1)

1 and port number with the IP address and NAT-generated port number  
2 of forwarder 302 and (ii) overwriting the destination IP address  
3 (VIP) with the IP address (H1) of the host (108(1)).

4 2. Forwarder 302 performs “Half-NAT” by overwriting  
5 the destination IP address (VIP) with the IP address (H1) of the host  
6 (108(1)) so that the source (client 102(1)) IP address (C1) and port  
7 number are preserved.

8 3. Forwarder 302 “tunnels” the packets received from  
9 client 102(1) from forwarder 302 to host 108(1). Specifically in this  
10 example, tunneling can be effectuated by encapsulating each packet  
11 within a new IP packet that is addressed to host 108(1). Network-  
12 load-balancing-aware software on host 108(1) reconstructs the  
13 original packet as received at forwarder 302 from client 102(1).  
14 This original packet is then indicated up on a virtual interface at host  
15 108(1) (e.g., the VIP address corresponding to forwarder 302 is  
16 bound to this virtual interface at host 108(1)). Exemplary  
17 implementations of such tunneling are described further below with  
18 reference to tunneler 312, especially for connection migration  
19 scenarios and particularly in the section entitled “Exemplary  
20 Connection Migrating with Optional Tunneling and/or Application-  
21 Level Load Balancing”.

22 Although FIGS. 4-9B show two specific separated functions, namely  
23 classifying and forwarding, it should be understood that other functions, such as  
24 those of request router 306, session tracker 308, connection migrator 310, and  
25 health and load handler 314, may also be scaled out independently (e.g., factored

1 out independently), as is described further below. Furthermore, it should be noted  
2 that one or more than two functions may be separated and scaled out  
3 independently at different times and/or simultaneously. Also, although TCP/IP is  
4 used for the sake of clarity in many examples in this and other sections, the  
5 network load balancing principles described herein are applicable to other  
6 transmission and/or communication protocols.

7 In the exemplary manner of FIG. 4, network load balancing functions (such  
8 as those shown in FIG. 3) may be separated from each other for scalability  
9 purposes. They may also be separated and duplicated into various configurations  
10 for increased availability. Exemplary configurations for scalability and/or  
11 availability are described below with reference to FIGS. 6-9B after the method of  
12 FIG. 5 is described.

13 FIG. 5 is a flow diagram 500 that illustrates an exemplary method for  
14 scaling out network load balancing infrastructure into different configurations.  
15 Flow diagram 500 includes three blocks 502-506. Although the actions of flow  
16 diagram 500 may be performed in other environments and with a variety of  
17 software schemes, FIGS. 1-4 and 6-9B are used in particular to illustrate certain  
18 aspects and examples of the method.

19 At block 502, network load balancing infrastructure is operated in a first  
20 configuration. For example, each configuration may relate to one or more of a  
21 selection, proportion, and/or interrelationship of different load balancing  
22 functionalities; a number of and/or type(s) of different devices; an organization  
23 and/or layout of different components; a distribution and/or allocation of  
24 resources; and so forth. At block 504, the network load balancing infrastructure is  
25 scaled out. For example, separated load balancing functionalities may be

1 expanded and/or concomitantly contracted on an individual and/or independent  
2 basis. At block 506, the scaled out network load balancing infrastructure is  
3 operated in a second configuration.

4 As noted above, a monolithic network load balancer may be scaled up by  
5 increasing network load balancing functionality in its entirety by supplanting  
6 previous network load balancing hardware with more-powerful network load  
7 balancing hardware. In contradistinction, scaling out network load balancing  
8 infrastructure can enable network load balancing (sub-)functions to be scaled out  
9 individually and/or independently. It can also enable network load balancing  
10 functions to be scaled out together or individually between and among different  
11 numbers of devices. Device, component, and resource-oriented scaling out  
12 examples are provided below.

13 FIG. 6 illustrates a first exemplary network load balancing infrastructure  
14 configuration from a device perspective. In this first device-oriented network load  
15 balancing infrastructure configuration, three devices 602(1), 602(2), and 602(3)  
16 are illustrated. However, one, two, or more than three devices 602 may  
17 alternatively be employed.

18 As illustrated, a forwarder 302(1), a classifier 304(1), and a host 108(1) are  
19 resident at and executing on device 602(1). A forwarder 302(2), a classifier  
20 304(2), and a host 108(2) are resident at and executing on device 602(2). Also, a  
21 forwarder 302(3), a classifier 304(3), and a host 108(3) are resident at and  
22 executing on device 602(3). Thus, in this first device-oriented network load  
23 balancing infrastructure configuration, a respective forwarder 302, classifier 304,  
24 and host 108 are sharing the resources of each respective device 602.  
25

1 In operation, forwarders 302 are the network endpoints for the VIP  
2 address(es). Any classifier 304 may plumb a route for a connection to any host  
3 108, depending on host status information. For example, classifier 304(2) may  
4 plumb a route for a new incoming connection to host 108(3). In accordance with a  
5 new route entry for this connection, forwarder 302(2) forwards subsequent packets  
6 to host 108(3).

7 In one alternative device-oriented network load balancing infrastructure  
8 configuration to which the illustrated first one may be scaled out, a fourth device  
9 602(4) (not explicitly shown in FIG. 6) may be added that includes a forwarder  
10 302(4), a classifier 304(4), and a host 108(4). If, on the other hand, sufficient  
11 classification functionality is already present with classifiers 304(1-3) but  
12 additional forwarding functionality can benefit the request handling of hosts 108, a  
13 fourth device 602(4) may be added that includes a forwarder 302(4) and optionally  
14 a host 108(4). For this scaled-out configuration, another classifier 304(1, 2, or 3)  
15 may plumb routes for forwarder 302(4) to any of hosts 108(1, 2, or 3) and host  
16 108(4), if present.

17 The first device-oriented exemplary network load balancing infrastructure  
18 configuration of FIG. 6 may be especially appropriate for smaller hosting  
19 situations in which separate devices for the network load balancing infrastructure  
20 are not technically and/or economically worthwhile or viable. However, as the  
21 hosting duties expand to a greater number (and/or a greater demand on the same  
22 number) of hosts 108 or if the network load on hosts 108 is significant, the first  
23 device-oriented exemplary network load balancing infrastructure configuration  
24 may be scaled out to accommodate this expansion, as represented by a second  
25

1 device-oriented exemplary network load balancing infrastructure configuration of  
2 FIG. 7.

3 FIG. 7 illustrates a second exemplary network load balancing infrastructure  
4 configuration from a device perspective. In this second device-oriented network  
5 load balancing infrastructure configuration, three devices 602(1), 602(2), and  
6 602(3) are also illustrated. Again, one, two, or more than three devices 602 may  
7 alternatively be employed.

8 As illustrated, forwarder 302(1) and classifier 304(1) are resident at and  
9 executing on device 602(1). Forwarder 302(2) and classifier 304(2) are resident at  
10 and executing on device 602(2). Also, forwarder 302(3) and classifier 304(3) are  
11 resident at and executing on device 602(3). Thus, in this second device-oriented  
12 network load balancing infrastructure configuration, each respective forwarder  
13 302 and classifier 304 are not sharing the resources of each respective device 602  
14 with a host 108. Furthermore, the network load balancing infrastructure may be  
15 servicing any number of hosts 108.

16 In operation, forwarders 302 are again the network endpoints for the VIP  
17 address(es). Also, any classifier 304 may plumb a route for a connection to any  
18 host 108, depending on host status information. For example, classifier 304(3)  
19 may plumb a route for a new incoming connection to host 108(2). In accordance  
20 with a new route entry for this connection, forwarder 302(3) forwards subsequent  
21 packets to host 108(2).

22 Hence, network load balancing infrastructure as realized in software, for  
23 example, may be scaled out by moving the network load balancing infrastructure  
24 (or part thereof) from devices that are shared with hosts 108 to devices that are not  
25 shared with hosts 108. Also, as alluded to above for FIG. 6, another device 602(4)



1 may be added to the network load balancing infrastructure to provide additional  
2 forwarding functionality, additional classifying functionality, additional  
3 functionality of both types, and so forth.

4 FIGS. 8A and 8B illustrate first and second exemplary network load  
5 balancing infrastructure configurations from a component perspective. As  
6 illustrated, first component-oriented exemplary network load balancing  
7 infrastructure configuration 800 includes four components. Second component-  
8 oriented exemplary network load balancing infrastructure configuration 850  
9 includes six components. An alternative second configuration 850 includes a  
10 seventh component as indicated by the dashed-line block, which is described  
11 further below.

12 Specifically, first component-oriented exemplary network load balancing  
13 infrastructure configuration 800 (or first configuration 800) includes (i) two  
14 forwarders 302(1) and 302(2) and (ii) two classifiers 304(1) and 304(2). Second  
15 exemplary component-oriented network load balancing infrastructure  
16 configuration 850 (or second configuration 850) includes (i) four forwarders  
17 302(1), 302(2), 302(3), and 302(4) and (ii) two classifiers 304(1) and 304(2).  
18 Thus, first configuration 800 is scaled out to second configuration 850 by adding  
19 two components, which are forwarding components in this example.

20 In a described implementation, each respective network-load-balancing-  
21 related functional component corresponds to a respective device (not explicitly  
22 shown in FIG. 8A or 8B); however, each component may alternatively correspond  
23 to part of a device or more than one device. For example, forwarders 302(1) and  
24 302(2) may be distributed across three devices. Or forwarder 302(1) and classifier  
25

1 304(1) may correspond to a first device, and forwarder 302(2) and classifier  
2 304(2) may correspond to a second device.

3 Two network-load-balancing-related functional components are added to  
4 scale out first configuration 800 to second configuration 850. However, one  
5 component (or more than two) may alternatively be added to scale out the network  
6 load balancing infrastructure. Furthermore, two or more different types of  
7 functional components may be scaled out “simultaneously”. For example, as  
8 illustrated by the dashed-line block, another classifying component (e.g., classifier  
9 304(3)) may also be added when scaling out first configuration 800 to second  
10 configuration 850.

11 Moreover, scaling by two or more different types of functional components  
12 may be performed in similar (e.g., equivalent) or dissimilar proportions to each  
13 other. As illustrated, adding forwarder components 302(3) and 302(4) while not  
14 adding any classifier component 304 or while adding a single classifier component  
15 304(3) represent a scaling out at dissimilar proportions. However, two classifier  
16 components 304(3) and 304(4) (the latter of which is not explicitly illustrated in  
17 FIG. 8B) may be added while the two forwarder components 302(3) and 302(4)  
18 are added for a scaling out at similar proportions. Regardless, each individual  
19 network-load-balancing-related functional component may consume a different  
20 amount of the available network load balancing infrastructure resources, as is  
21 described with reference to FIGS. 9A and 9B.

22 FIGS. 9A and 9B illustrate first and second exemplary network load  
23 balancing infrastructure configurations from a resource perspective. First  
24 resource-oriented exemplary network load balancing infrastructure configuration  
25 900 (or first configuration 900) includes a first resource distribution or allocation

1 for a load balancing unit 106. Second resource-oriented exemplary network load  
2 balancing infrastructure configuration 950 (or second configuration 950) includes  
3 a second resource distribution for load balancing unit 106.

4 As illustrated, first configuration 900 includes a 70%-30% resource  
5 distribution, and second configuration 950 includes a 40%-60% resource  
6 distribution. Such resources may include total device resources (e.g., number of  
7 devices), processing resources (e.g., number of processor cycles), memory  
8 resources (e.g., portion of cache, main memory, etc.), network bandwidth and/or  
9 interface resources (e.g., bits per second and/or physical network interface cards  
10 (NICs)), and so forth.

11 Specifically for first configuration 900, forwarder 302 consumes 70% of  
12 the resources of load balancing unit 106 while classifier 304 consumes 30% of  
13 these resources. After reallocation during a scaling out procedure to produce  
14 second configuration 950, forwarder 302 consumes 40% of the resources of load  
15 balancing unit 106 while classifier 304 consumes 60% of these resources.

16 In an exemplary situation, first configuration 900 might facilitate better  
17 network load balancing performance when fewer, longer transactions are being  
18 handled by the associated hosts (not shown in FIGS. 9A and 9B) because  
19 classification functionality is utilized upon initial communication for a connection  
20 and forwarding functionality is utilized thereafter. Second configuration 950, on  
21 the other hand, might facilitate better network load balancing performance when  
22 more, shorter transactions are being handled by the associated hosts because the  
23 classification functionality is utilized for a greater percentage of the total number  
24 of packets funneled through the network load balancing infrastructure. In this  
25 situation, if request routing functionality is also being employed, then request

1 router(s) 306 are also allocated a percentage of the total computing resources. The  
2 resource distribution among the three functionalities may be adjusted while  
3 handling connections (e.g., adjusted “on the fly”) depending on current resource  
4 consumption and/or deficits.

5 As indicated above with reference to FIGS. 2 and 3, each load balancing  
6 unit 106 may correspond to all or a part of a total network load balancing  
7 infrastructure 106. For any given physically, logically, arbitrarily, etc. defined or  
8 stipulated load balancing unit 106, the resources thereof may be re-allocated  
9 during a scale out procedure. More specifically, a resource distribution  
10 between/among different network-load-balancing-related separated functions of a  
11 load balancing unit 106 may be altered in a scale out procedure. Furthermore,  
12 more than two different functions, as well as other network-load-balancing-related  
13 functions that are not specifically illustrated in FIGS. 9A and 9B, may be allocated  
14 differing resource percentages.

15 The percentage of total system resources allocated to all load balancing  
16 functions may also be altered in a scale out procedure. As a general processing  
17 power example, the percentage of total processing power that is devoted to load  
18 balancing may be gradually increased as the amount of traffic that needs to be load  
19 balanced increases.

20 Network load balancing software may optionally perform monitoring to  
21 analyze and determine whether resources should be reallocated. For example, the  
22 network load balancing software may monitor the processor utilization of different  
23 network-load-balancing-related functions. The actual reallocation may also  
24 optionally be performed automatically by the network load balancing software in  
25 an offline or online mode.

1 It should be understood that a scaling out capability of network load  
2 balancing infrastructure (e.g., as realized at least partially in software) as described  
3 herein may relate to different installations and not necessarily a change to a single  
4 installation. In a resource-oriented example, network load balancing infrastructure  
5 as described herein may be configured in accordance with one resource  
6 distribution in one installation environment and may be configured in accordance  
7 with another different resource distribution in another installation environment  
8 having different operational parameters. Additionally, the capabilities, features,  
9 options, etc. described above with regard to scaling out are also applicable for  
10 “scaling in”. In other words, resources devoted to network load balancing  
11 infrastructure (or sub-functions thereof) may also be reduced.

#### 12 **Exemplary Health and Load Handling**

13 This section describes how host status information, such as health and/or  
14 load information, may be collected for and utilized in network load balancing.  
15 This section primarily references FIGS. 10-18 and illuminates health and load  
16 functionality such as that provided by health and load handler 314 (of FIG. 3). As  
17 described above with reference to FIG. 3, each host 108 hosts one or more  
18 applications 316. Health and load handler 314 utilizes health and/or load  
19 information that relates to applications 316 and/or hosts 108 for certain described  
20 implementations of network load balancing.

21 FIG. 10 illustrates an exemplary network load balancing approach that  
22 involves host status information (HSI) 1006. Each host 108(1), 108(2) ... 108(n)  
23 includes one or more applications 316(1), 316(2) ... 316(n), respectively. These  
24 hosts 108 generally and these applications 316 specifically may change statuses  
25 from time to time.

1 For example, hosts 108 and applications 316 may be accepting new  
2 connections or not accepting new connections. Also, they may be quickly  
3 handling client requests or slowly handling client requests. Furthermore, they may  
4 have many resources in reserve or few unused resources. All or any part of such  
5 data, or other data, may comprise host status information 1006. Generally, host  
6 status information 1006 provides an indication of the status of some aspect of  
7 hosts 108 and/or applications 316 that are running thereon.

8 In a described implementation, each host 108(1), 108(2) ... 108(n) includes  
9 a host status information (HSI) determiner 1002(1), 1002(2) ... and 1002(n),  
10 respectively. Each host 108(1), 108(2) ... 108(n) also includes a host status  
11 information (HSI) disseminator 1004(1), 1004(2) ... and 1004(n), respectively.  
12 Each host status information determiner 1002 and/or host status information  
13 disseminator 1004 may be part of load balancing infrastructure (LBI) 106.

14 Each host status information determiner 1002 determines host status  
15 information 1006 for its respective host 108 and/or applications 316 that are  
16 running thereon. Exemplary techniques for determining such host status  
17 information 1006 are described below with reference to FIGS. 12-14, and  
18 particularly FIG. 13A. Each host status information disseminator 1004  
19 disseminates host status information 1006 for its respective host 108 and/or  
20 applications 316 to load balancing infrastructure 106 (e.g., those portion(s) of load  
21 balancing infrastructure 106 that are not located on hosts 108). Exemplary  
22 techniques for disseminating such host status information 1006 are described  
23 below with reference to FIGS. 12-17, and particularly FIGS. 13B and 15-17.

24 Specifically, each host status information disseminator 1004 disseminates  
25 host status information 1006 (directly or indirectly) to each load balancing unit

1 (LBU) 106 of load balancing infrastructure 106 that includes at least one health  
2 and load handler 314 and/or classifier 304. Load balancing infrastructure 106  
3 refers to host status information 1006 when implementing network load balancing.  
4 For example, as indicated by logic 1008, load balancing infrastructure 106 is  
5 capable of making load balancing decisions responsive to host status information  
6 1006.

7 In operation at (1), host status information determiners 1002 determine host  
8 status information 1006 for respective hosts 108 and/or applications 316. At (1)  
9 and (2), host status information disseminators 1004 disseminate host status  
10 information 1006 from hosts 108 to load balancing infrastructure 106. For  
11 example, host status information 1006 may be disseminated to individual load  
12 balancing units 106. At (3), logic 1008 makes network load balancing decisions  
13 responsive to host status information 1006. At (4), connections are forwarded to  
14 targeted hosts 108 based on these network load balancing decisions.

15 FIG. 11 is a flow diagram 1100 that illustrates an exemplary method for  
16 network load balancing that involves host status information. Flow diagram 1100  
17 includes three blocks 1102-1106. Although the actions of flow diagram 1100 may  
18 be performed in other environments and with a variety of software schemes, FIGS.  
19 1-3 and 10 are used in particular to illustrate certain aspects and examples of the  
20 method.

21 At block 1102, host status information is sent from hosts to load balancing  
22 units. For example, host status information 1006 may be sent from hosts 108 to  
23 load balancing units 106. At block 1104, the host status information is received  
24 from the hosts at the load balancing units. For example, load balancing units 106  
25 may receive host status information 1006 from hosts 108. At block 1106, load

1 balancing decisions are made responsive to the received host status information.  
2 For example, logic 1008 at load balancing units 106 may make decisions for  
3 network load balancing responsive to host status information 1006.

4 Thus in FIG. 10, load balancing infrastructure 106 collects host status  
5 information 1006 from hosts 108 (and/or applications 316 thereof) and load  
6 balances incoming requests that are directed to hosts 108 responsive to host status  
7 information 1006. As described further below with reference to FIGS. 12-18, this  
8 host status information 1006 may be application-specific. As also described  
9 further below, examples of host status information 1006 include health and/or load  
10 information.

11 FIG. 12 illustrates an exemplary network load balancing approach that  
12 involves health and/or load information (HLI) 1206. Hosts 108(1), 108(2) ...  
13 108(n) are coupled to load balancing units 106(1), 106(2) ... 106(u) via a  
14 communication linkage 1210 such as a network.

15 As illustrated, hosts 108 communicate health and load information 1206 to  
16 load balancing units 106 using communication linkage 1210. The bi-directional  
17 communication of health and load information 1206, as indicated by the double-  
18 pointed arrow, refers to a two-way communication from load balancing units 106  
19 to hosts 108 that provides certain completeness, coherency, correctness, etc. such  
20 that hosts 108 and/or load balancing units 106 may fail independently of one  
21 another. Such two-way communications from load balancing units 106 to hosts  
22 108 are described further below with particular reference to FIG. 15.

23 Health information reflects whether a given host and/or application is  
24 capable of handling client requests. Load information reflects the number,  
25 amount, and/or level of client requests that the given host and/or application is



1 capable of handling at a particular moment. In other words, load can reflect  
2 directly and/or inversely an available number, amount, and/or level of total  
3 capacity of the given host and/or application. As noted above, implementations  
4 described with reference to FIGS. 12-18 focus on health and/or load information;  
5 however, those implementations are also applicable to general status information  
6 for hosts (including the applications thereof).

7 In a described implementation, each host 108(1), 108(2) ... 108(n) includes  
8 a respective health and load infrastructure (H&LI) component 1202(1), 1202(2) ...  
9 1202(n). Each health and load infrastructure component 1202 may optionally be a  
10 portion of load balancing infrastructure 106 that is resident at and executing on  
11 each host 108. Health and load information 1206 may be realized in software.  
12 When functioning, each health and load infrastructure 1202(1), 1202(2) ...  
13 1202(n) creates and maintains a respective health and load (H&L) table 1204(1),  
14 1204(2) ... 1204(n).

15 These health and load tables 1204 may include application-specific entries.  
16 Health and load information 1206 that is stored in health and load tables 1204 may  
17 be independent of load balancing infrastructure 106. For example, administrators,  
18 designers, etc. may specify criteria for health and load information 1206 at  
19 configuration time. Additionally, entities external to a device that is or that has a  
20 host 108 may contribute to determining health and load information 1206 for  
21 applications 316 on the device. An exemplary health and load table 1204 is  
22 described further below with reference to FIG. 13A.

23 Each load balancing unit 106(1), 106(2) ... 106(u) includes a respective  
24 consolidated health and load (H&L) cache 1208(1), 1208(2) ... 1208(u). Each  
25 consolidated health and load cache 1208 includes the information from each health

1 and load table 1204(1), 1204(2) ... 1204(n). Consequently, each load balancing  
2 unit 106 is provided with quick (e.g., cached) access to health and load  
3 information 1206 for each host 108 for which load balancing units 106 are load  
4 balancing network traffic.

5 In operation, health and load infrastructures 1202 push health and load  
6 information 1206 from health and load tables 1204 to consolidated health and load  
7 caches 1208. The mechanism to provide health and load information 1206 is  
8 event driven such that changes to health and load tables 1204 are provided to  
9 consolidated health and load caches 1208 in a timely, scaleable manner.

10 FIG. 13A is an exemplary health and load table 1204 as illustrated in FIG.  
11 12. In a described implementation, health and load table 1204 includes multiple  
12 entries 1302 that are each associated with a different application 316. Each entry  
13 1302 may correspond to a row in health and load table 1204 that has three  
14 columns. These columns correspond to application identifier (ID) 1302(A),  
15 application status characterization 1302(B), and load balancer directive 1302(C).

16 Because each entry 1302 is associated with a particular application 316, a  
17 row is added as each application is spun up (e.g., by an administrator). Likewise,  
18 a row is deleted/removed each time an application is closed down. Similarly,  
19 individual fields in columns 1302(A), 1302(B), and/or 1302(C) are  
20 modified/updated when a value thereof changes. For example, when a status  
21 characterization value changes for a given application 316, a value in a field of  
22 application status characterization 1302(B) for entry 1302 of the given application  
23 316 is updated.

24 The additions and deletions of entries 1302 for applications 316 may be  
25 effectuated with input from a control manager at the host 108. For example, a

1 control manager portion of an operating system knows when an application 316 is  
2 started and stopped because it is actively involved in the starting and stopping of  
3 applications 316. Hence, a control manager may identify that it has, at least in  
4 part, started an application 316, and the control manager may establish that it has,  
5 at least in part, stopped the application 316. Health and load infrastructure 1202  
6 may therefore be informed of the starting and stopping of applications 316 by the  
7 control manager. Hence, no such explicit communication from applications 316  
8 has to be provided to health and load infrastructure 1202. An example of a control  
9 manager is the Service Control Manager (SCM) of the Windows® Operating  
10 System from Microsoft® Corporation.

11 Application identifier 1302(A) includes information that is used to uniquely  
12 identify the application 316 to which entry 1302 is associated. Application  
13 identifier 1302(A) may include one or more of the following for the associated  
14 application 316: the virtual IP address and port, the physical IP address and port,  
15 the protocol used, and any protocol-specific information. The protocol may be  
16 HTTP, IPsec, SOAP, and so forth. The protocol-specific information may be a  
17 URL pattern or string to further delineate the application associated with entry  
18 1302. Thus, application identifier 1302(A) more particularly refers to a specific  
19 application endpoint on a particular host 108.

20 Other application identifiers may alternatively be employed. For example,  
21 to reduce communication bandwidth, application identifier 1302(A) may be a 32-  
22 bit number that maps to the above exemplary information at health and load  
23 infrastructure 1202 and at load balancing units 106. Moreover, any of the fields in  
24 entry 1302 may actually contain a globally unique identifier (GUID) that is used  
25 as a key to look up the true information for the field.

1       Application status characterization 1302(B) includes information that  
2 reflects the status of the application 316 to which entry 1302 is associated.  
3 Application status characterization 1302(B) includes the following for the  
4 associated application 316: application health, application load, and application  
5 capacity. Application health is a quasi-Boolean value that indicates whether an  
6 application is functioning. Application health can be healthy, failing, or unknown.  
7 Application health is a relatively-instantaneous value and is communicated with  
8 relatively low latency (e.g., of approximately a second or a few seconds) to load  
9 balancing units 106 when the application health value changes.

10       Application load is a value that indicates how occupied or busy a given  
11 application is and thus, directly or inversely, how much additional load the given  
12 application can handle. Application load is a relatively slowly-changing or  
13 averaged value that can be smoothed with a hysteresis-inducing mechanism, if  
14 desired, to eliminate transient spikes of increased or decreased load. It is  
15 communicated relatively infrequently to load balancing units 106 (e.g.,  
16 approximately one to four times a minute). The value of application load is given  
17 meaning with regard to application capacity.

18       Application capacity is a value that indicates the maximum capacity of the  
19 application. It is selected in a generic manner to be meaningful for a given context  
20 but still sufficiently flexible for other contexts. Application capacity is a unit-less,  
21 bounded number (e.g., 0-99) that is determinable at configuration time. It may be  
22 based on processing power, memory size/speed, network access, some  
23 combination thereof, and so forth. Application capacity expresses relative  
24 capacities between and among other applications of the same type in a set of hosts  
25 108(1, 2 ... n).

1        Thus, relative to application capacity, application load gains meaning.  
2        Application load for a given application is a percentage of application capacity for  
3        the given application. Alternatively, application load can be expressed as a unit-  
4        less number from which the percentage can be ascertained in conjunction with the  
5        value of application capacity.

6        Load balancer directive 1302(C) includes information that reflects the  
7        desired and/or expected state of the directive established by health and load  
8        infrastructure 1202 for load balancing units 106 with respect to an application 316  
9        to which entry 1302 is associated. Load balancer directive 1302(C) includes the  
10       following for the associated application 316: target load balancing state and  
11       current load balancing state.

12       The target load balancing state reflects the state of the directive to load  
13       balancing units 106 as desired by health and load infrastructure 1202. The current  
14       load balancing state reflects what health and load infrastructure 1202 understands  
15       the current state of the directive to load balancing units 106 to be as recorded at  
16       load balancing units 106. The current load balancing state thus reflects the load  
17       balancing directive that health and load infrastructure 1202 expects load balancing  
18       units 106 to be currently operating under as dictated using a communication  
19       protocol. Such an exemplary communication protocol is described further below  
20       with reference to FIG. 15. The interaction and relationship between the target load  
21       balancing state and the current load balancing state is also further clarified with  
22       the description of FIG. 15.

23       The target load balancing state and the current load balancing state may  
24       each take a value of active, inactive, or draining. An active directive indicates that  
25       new requests/connections are welcome and may be targeted at the application that

1 is associated with entry 1302. An inactive directive indicates that no additional  
2 packets should be forwarded to the associated application. A draining directive  
3 indicates that no packets for new requests/connections should be sent to the  
4 associated application but that packets for existing requests/connections should  
5 continue to be forwarded to the associated application.

6 In a described implementation, the definitive version of respective health  
7 and load information 1206 is stored at health and load tables 1204 that are located  
8 at each respective host 108 of multiple hosts 108. With this implementation, if a  
9 host 108 crashes, the health and load information 1206 that is lost pertains to those  
10 applications 316 that also crashed. A measure of high availability is therefore  
11 gained automatically without duplicating data. However, the definitive version of  
12 health and load information 1206 may alternatively be stored elsewhere. Other  
13 such storage options include load balancing units 106 themselves, a host 108 that  
14 (as its sole task or along with hosting duties) stores and maintains health and load  
15 information 1206 for multiple other (including all other) hosts 108, another  
16 separate and/or external device, and so forth.

17 If the definitive version of health and load information 1206 is stored and  
18 maintained elsewhere besides being distributed across hosts 108(1, 2 ...n), such  
19 health and load information 1206 may be stored redundantly (e.g., also stored in a  
20 duplicative device, backed-up, etc.) for high-availability purposes. Exemplary  
21 proxy scenarios for storing health and load information 1206 are described below  
22 with reference to FIGS. 17A and 17B. FIG. 17A is directed to a proxy scenario for  
23 health and load tables 1204, and FIG. 17B is directed to a proxy scenario for  
24 consolidated health and load caches 1208.  
25

1        FIG. 13B is an exemplary consolidated health and load cache 1208 as  
2        illustrated in FIG. 12. In a described implementation, each consolidated health and  
3        load cache 1208 in each load balancing unit 106 includes at least part of the  
4        information stored in each health and load table 1204 for each health and load  
5        infrastructure 1202 at each host 108. The cached health and load information may  
6        be organized in any manner in consolidated health and load cache 1208.

7        As illustrated, consolidated health and load cache 1208 includes a cache for  
8        each host 108(1), 108(2) ... 108(n) that replicates part or all of the information in  
9        the health and load table 1204 of each respective host 108(1, 2 ... n). Specifically,  
10       consolidated health and load cache 1208 includes a cache for host #1 1304(1), a  
11       cache for host #2 1304(2) ... a cache for host #n 1304(n). Thus, the illustrated  
12       consolidated health and load cache 1208 is organized at a broad level by host  
13       108(1, 2 ... n), with each individual cache 1304 including application-specific  
14       entries for the corresponding respective host 108(1, 2 ... n). Alternatively,  
15       consolidated health and load cache 1208 may be organized at a broad level by type  
16       of application 316, with individual blocks that are directed to a specific application  
17       type further divided by host 108(1, 2 ... n). Other data structure formats may also  
18       be employed.

19       FIG. 14 is a flow diagram that illustrates an exemplary method for network  
20       load balancing that involves health and load information. Flow diagram 1400  
21       includes eight blocks 1402-1416. Although the actions of flow diagram 1400 may  
22       be performed in other environments and with a variety of software schemes, FIGS.  
23       1-3 and 12-13B are used in particular to illustrate certain aspects and examples of  
24       the method. For example, the actions of two blocks 1402-1404 are performed by a  
25

1 host 108, and the actions of six blocks 1406-1416 are performed by a load  
2 balancing unit 106.

3 At block 1402, health and load information at a host is determined. For  
4 example, health and load information 1206 for applications 316(2) may be  
5 ascertained by health and load infrastructure 1202(2) and stored in health and load  
6 table 1204(2) at host 108(2). At block 1404, the determined health and load  
7 information is disseminated to load balancing units. For example, health and load  
8 infrastructure 1202(2) may send health and load information 1206 for applications  
9 316(2) to load balancing units 106(1, 2 ... u). As indicated by arrow 1418, the  
10 actions of blocks 1402 and 1404 are repeated so that (application) health and load  
11 may be continually monitored and updated as changes occur.

12 At block 1406, health and load information is received from hosts. For  
13 example, load balancing unit 106(1) may receive health and load information 1206  
14 from multiple hosts 108(1, 2 ... n), which includes health and load information  
15 1206 for applications 316(2) of host 108(2). At block 1408, the received health  
16 and load information is cached. For example, load balancing unit 106(1) may  
17 store health and load information 1206 from hosts 108(1, 2 ... n) into consolidated  
18 health and load cache 1208(1). With reference to the FIG. 13B implementation of  
19 a consolidated health and load cache 1208(1), health and load information 1206  
20 for applications 316(2) from host 108(2) may be stored in cache for host #2  
21 1304(2). As indicated by arrow 1420, the actions of blocks 1406 and 1408 are  
22 repeated so that (application) health and load information may be continually  
23 received and updated as changes occur.

24 As indicated by dashed arrow 1422, load balancing units 106 are also  
25 handling communications from clients 102 while handling (application) health and



1 load issues. At block 1410, a packet requesting a new connection is received. For  
2 example, load balancing unit 106(1) may receive a TCP SYN packet from client  
3 102(2) through network 104. At block 1412, the cached health and load  
4 information is consulted. For example, load balancing unit 106(1) may consult  
5 consolidated health and load cache 1208(1). More particularly, load balancing unit  
6 106(1) may consult entries that are associated with the application to which the  
7 TCP SYN packet is directed across caches for hosts #1, #2 ... #n 1304(1, 2 ... n).

8 At block 1414, a host is selected responsive to the cached health and load  
9 information. For example, load balancing unit 106(1) may select host 108(2)  
10 having application(s) 316(2) responsive to health and load information 1206 that is  
11 cached in consolidated health and load cache 1208(1). The selected application  
12 316 (and host 108) should be healthy and able to accept additional load (e.g.,  
13 possibly the least loaded application among those applications that are of the  
14 application type to which the TCP SYN packet is directed).

15 The consulting of the cached health and load information (at block 1412)  
16 and the host-selecting responsive to the cached health and load information (at  
17 block 1414) may be performed prior to reception of a specific new-connection-  
18 requesting packet and/or using a batched scheme. Also, the selecting may be in  
19 accordance with any of many schemes. For example, a token based or a round-  
20 robin based scheme may be employed. With either scheme, the selection may  
21 involve a weighting of relative loads among the application options. This  
22 consultation and selection, along with the token and round-robin based schemes,  
23 are described further below with reference to FIG. 18 and in the section entitled  
24 "Exemplary Classifying, Forwarding, and Request Routing", especially with  
25 regard to classifying functionality.

1 After the target host is selected at block 1414, the new-connection-  
2 requesting packet may be sent thereto. At block 1416, the packet received from  
3 the client is forwarded to the selected host. For example, the TCP SYN packet is  
4 forwarded from load balancing unit 106(1) to selected host 108(2). The  
5 forwarding of this initial packet may be effectuated directly by a classifier 304 or  
6 by a forwarder 302, as is also described further below in the section entitled  
7 "Exemplary Classifying, Forwarding, and Request Routing".

8 For a described implementation, health and load infrastructure 1202 is  
9 resident at and distributed across multiple hosts 108 as well as being located at  
10 load balancing units 106 (as represented by health and load handler 314). Health  
11 and load infrastructure 1202 has three responsibilities. First, it exposes listening  
12 point(s) to attain application status updates for application status characterizations  
13 1302(B) of health and load tables 1204. Second, it synthesizes the application  
14 status information to determine what load balancing units 106 should do, which is  
15 embodied in load balancer directive 1302(C). Third, health and load infrastructure  
16 1202 communicates this directive from hosts 108 to load balancing units 106.

17 The directive content of load balancer directive 1302(C) is effectively a  
18 digested version of the information for application status characterizations  
19 1302(B). However, load balancing units 106 may also receive the raw information  
20 of application status characterizations 1302(B) as well as this processed directive.  
21 The communication of the content of these and other fields of health and load  
22 tables 1204 is accomplished using a message protocol that is described below with  
23 reference to FIG. 15.

24 FIG. 15 illustrates an exemplary message protocol 1500 for the health and  
25 load information-related communications that are illustrated in FIG. 12 between

1 hosts 108 and load balancing units 106. Generally, an event-driven mechanism is  
2 used to push changes to health and load tables 1204 from hosts 108 to load  
3 balancing units 106. In other words, for a described implementation, information  
4 is transmitted from hosts 108 to load balancing units 106 when health and load  
5 tables 1204 are updated. This avoids periodically sending a snapshot of all of each  
6 health and load table 1204, which reduces network bandwidth consumption by  
7 health and load infrastructure 1202.

8 Message protocol 1500 may be implemented using any available message  
9 transport mechanism. Such mechanisms include reliable multicast transmission,  
10 point-to-point transmission (e.g., user datagram protocol (UDP)), and so forth. As  
11 illustrated, message protocol 1500 includes seven message types 1502-1514: a  
12 heartbeat message 1502, a goodbye message 1504, a row change message 1506, a  
13 get table snapshot message 1508, a send table snapshot message 1510, a postulate  
14 table state message 1512, and a postulate wrong message 1514.

15 It should be understood that, with the exception of arrows 1516 and 1518,  
16 no temporal relationship between or among the different messages types 1502-  
17 1514 is implied by the illustration. For example, a row change message 1506 does  
18 not typically follow a goodbye message 1504.

19 Heartbeat message 1502 indicates that a particular host 108 is functioning  
20 and provides some error checking for the content of a corresponding particular  
21 health and load table 1204 with respect to a corresponding particular cache for the  
22 particular host 1304 in consolidated health and load cache 1208. Each health and  
23 load infrastructure 1202 at each host 108 sends a heartbeat message directly or  
24 indirectly to each consolidated health and load cache 1208 at each load balancing  
25 unit 106.

1 Heartbeat messages 1502 address the aging-out problem for data in  
2 consolidated health and load caches 1208 that arises, in part, because a snapshot of  
3 the entirety of each health and load table 1204 is not periodically transmitted to  
4 each load balancing unit 106. A transmission scheme for heartbeat messages 1502  
5 is described further below with reference to FIG. 16.

6 Heartbeat messages 1502 include an identifier for the host, error checking  
7 data, and optionally a DNS name. The identifier of the host may be a unique (e.g.,  
8 32-bit) number that is selected at configuration time. The error checking data may  
9 be a checksum, a state-change sequence number, a generation number, a CRC  
10 value, etc. that enables a receiving load balancing unit 106 to validate that the  
11 contents of its consolidated health and load cache 1208 comports with the contents  
12 of the health and load table 1204 of the transmitting host 108. If a generation  
13 number approach is employed, then multiple generation IDs can be used with each  
14 generation ID assigned to a “chunk” of applications. Messages can then refer to a  
15 chunk number or a chunk number/generation ID pair, depending on the context.

16 The error checking data (or, more generally, a content indicator) may be a  
17 single value for the health and load table 1204 overall, or it may be multiple values  
18 determined on a per-entry 1302 basis. The DNS name may optionally be sent  
19 (e.g., every “x” heartbeats) to verify or update the current correct network address  
20 for the host.

21 Goodbye message 1504 is sent from a particular host 108 to load balancing  
22 units 106 to indicate that the particular host 108 is planning to shutdown.  
23 Goodbye message 1504 includes a host identifier that may be indexed/mapped to a  
24 network address for the particular host 108. Goodbye message 1504 is used for  
25 clean, intentional shutdowns by hosts 108 to precipitate a “fast clear”. However, if

1 a goodbye message 1504 is lost, caches eventually age out the particular host's  
2 108 entries because heartbeat messages 1502 are no longer sent.

3 Row change message 1506 is sent from a particular host 108 to load  
4 balancing units 106 to indicate that the health and/or load for a given application  
5 316 of the particular host 108 has changed. Row change message 1506 includes a  
6 host identifier, an application identifier, an operation, and data for the operation.  
7 Exemplary host identifiers are described above with regard to heartbeat messages  
8 1502 and goodbye messages 1504. Exemplary application identifiers are  
9 described above with regard to application identifier 1302(A) of an application-  
10 associated entry 1302 of health and load tables 1204.

11 The row change operation may be add, delete, or update. In other words,  
12 the data for the operation may be added to (for an add operation) or a replacement  
13 for (for an update operation) information already present at consolidated health  
14 and load caches 1208 at load balancing units 106. For a delete operation, no data  
15 need be provided. Message protocol 1500 is defined such that multiple operations  
16 may be stipulated to be performed for a single row change message 1506. Hence  
17 for a particular host identifier, sets of an application identifier, operation, and  
18 operation data may be repeated for multiple applications 316 of the host 108 that is  
19 identified by the particular host identifier.

20 Get table snapshot message 1508 is sent from a particular load balancing  
21 unit 106 for a particular consolidated health and load cache 1208 to an individual  
22 host 108 or hosts 108. This get table snapshot message 1508 requests that health  
23 and load infrastructure 1202 at hosts 108 provide a snapshot of the respective  
24 health and load table 1204 for the respective host 108. This message includes an  
25 identification of the requesting load balancing unit 106 and may be used by a load

balancing unit 106 (i) after it has failed and then recovered; (ii) after a host 108 fails, recovers, and begins sending heartbeat messages 1502 again; (iii) if a row change message 1506 is sent to load balancing unit 106, but the message gets dropped, so its consolidated health and load cache 1208 is out of sync with the respective health and load table 1204 for the respective host 108; and (iv) so forth.

For the third (iii) situation, the lack of synchronization between consolidated health and load cache 1208 and the respective health and load table 1204 for the respective host 108 is discovered by a subsequent heartbeat message 1502 from the respective host 108 because the “error checking” will indicate that consolidated health and load cache 1208 is out of date. Load balancing unit 106 can then send a get table snapshot message 1508 so that it can update its consolidated health and load cache 1208. Thus, for any of the three (i, ii, iii) exemplary situations, load balancing unit 106 subsequently reconstitutes its consolidated health and load cache 1208 using get table snapshot 1508. Get table snapshot 1508 may be sent repeatedly to each host 108 in a point-to-point manner or may be sent one time to many hosts 108 in a multicast manner.

Send table snapshot message 1510 is sent from an individual host 108 to a particular load balancing unit 106 after the individual host 108 has received a get table snapshot message 1508 from the particular load balancing unit 106 as indicated by arrow 1516. The contents of a send table snapshot message 1510 is prepared by health and load infrastructure 1202 and may include all or at least multiple rows of the health and load table 1204 of the individual host 108 so that the particular load balancing unit 106 may rebuild its consolidated health and load cache 1208. Send table snapshot message 1510 may be a separately designed

1 message, or it may be equivalent to a sequence of add operations in a row change  
2 message 1506.

3 Postulate table state message 1512 and postulate wrong message 1514 are  
4 related to the target load balancing state and the current load balancing state of  
5 load balancer directive 1302(C) of an entry 1302 in a health and load table 1204.  
6 The target load balancing state is the directive that health and load infrastructure  
7 1202 desires load balancing units 106 to be operating under. The current load  
8 balancing state is the directive that health and load infrastructure 1202 expects or  
9 believes that load balancing units 106 are currently operating under. Generally,  
10 the two load balancing states are identical.

11 However, the target load balancing state may differ from the current load  
12 balancing state during a transitional period for a state directive change. For  
13 example, the target load balancing state and the current load balancing state are  
14 both initially set to active. A problem with host 108 and/or an application 316  
15 thereof is detected and the target load balancing state directive is switched to  
16 draining. This draining directive is communicated to load balancing units 106  
17 using a row change message 1506.

18 There is a delay before this directive change is noted in all consolidated  
19 health and load caches 1208 of all load balancing units 106. During this  
20 transitional period, the target load balancing state is draining while the current load  
21 balancing state is still active at health and load table 1204 of host 108. Before  
22 changing the current load balancing state to draining, health and load  
23 infrastructure 1202 wants to ensure that consolidated health and load caches 1208  
24 have actually been updated to the new directive state of draining.  
25

1 To verify that consolidated health and load caches 1208 of load balancing  
2 units 106 have been updated to a new state directive, health and load infrastructure  
3 1202 sends a postulate table state message 1512 to load balancing units 106.  
4 Postulate table state message 1512 is sent some time (e.g., a predetermined delay  
5 period) after transmission of a row change message 1506 indicating that the state  
6 directive is to be changed. The postulate table state message 1512, in this  
7 example, indicates that the table state should be draining. As indicated by the  
8 dashed arrow 1518, a load balancing unit 106 responds to this postulate table state  
9 message 1512 if its consolidated health and load cache 1208 differs from the  
10 postulated state directive.

11 If the directive in consolidated health and load cache 1208 does differ from  
12 the postulated state directive, then that load balancing unit 106 sends a postulate  
13 wrong message 1514 to the health and load infrastructure 1202 of the host 108 that  
14 issued the postulate table state message 1512. This health and load infrastructure  
15 1202 then periodically resends postulate table state message 1512 until no further  
16 postulate wrong messages 1514 are received from consolidated health and load  
17 caches 1208. At that point, health and load infrastructure 1202 sends a row change  
18 message 1506 with the new current load balancing state. In this sense,  
19 consolidated health and load caches 1208 are the definitive determiners of the  
20 current load balancing state, and health and load infrastructure 1202 is the  
21 definitive determiner of the target load balancing state.

22 FIG. 16 illustrates an exemplary message transmission scheme for the  
23 communications that are illustrated in FIG. 12 between hosts 108 and load  
24 balancing units 106. The exemplary message transmission scheme can reduce the  
25 bandwidth consumed by heartbeat messages 1502 on communication linkage



1 1210. The message transmission scheme of FIG. 16 is particularly adapted to  
2 heartbeat messages 1502, but it may also be utilized for other messages of  
3 message protocol 1500.

4 A group of hosts 108(1), 108(2), 108(3) ... 108(11), and 108(12) are  
5 illustrated along with load balancing units 106(1), 106(2) ... 106(u). Each line  
6 represents membership linkage or inclusion among the group of hosts 108(1, 2 ...  
7 12). The group of hosts 108(1, 2 ... 12) form a membership of nodes that work  
8 together to propagate heartbeat information to load balancing units 106. Although  
9 twelve hosts are shown, more or fewer may be part of any given group of hosts.  
10 Also, a total set of hosts 108 that are being served by a load balancing  
11 infrastructure 106 may be divided into one, two, three, or more groups of hosts.

12 In a described implementation, the membership of nodes for group of hosts  
13 108(1, 2 ... 12) elect a leader that is responsible for transmitting heartbeat  
14 messages 1502 to load balancing units 106. Each (non-leading) host 108 in group  
15 of hosts 108(1, 2 ... 12) sends its heartbeat messages 1502 to the elected leader.  
16 Host 108(4) is the elected leader in this example.

17 With the membership of nodes, heartbeat information for each host 108 in  
18 group of hosts 108(1, 2 ... 12) propagates to the group leader host 108(4). Host  
19 108(4) collects the heartbeat information and consolidates it into a consolidated  
20 heartbeat message 1602. Consolidated heartbeat messages 1602(1), 1602(2) ...  
21 1602(u) are then sent to respective load balancing units 106(1), 106(2) ... 106(u).  
22 These consolidated heartbeat messages 1602 may optionally be compressed to  
23 further reduce bandwidth consumption.

24 As another alternative, the leader host 108(4) may only forward changes in  
25 group membership to consolidated health and load caches 1208. In other words,

1 in this mode, consolidated health and load caches 1208 deal primarily if not solely  
2 with state changes to membership. It is the responsibility of the leader host 108(4)  
3 to ensure that the first hello is forwarded when a host 108 comes online and that a  
4 goodbye message 1504 gets sent when that host 108 goes offline. Additionally, a  
5 host 108 can periodically specify that a heartbeat message 1502 is to be  
6 "forwarded". This indicates to the leader host 108(4) to send it to consolidated  
7 health and load caches 1208 even though it does not represent a membership  
8 change.

9 Heartbeat messages 1502 (including consolidated heartbeat messages 1602)  
10 are used by load balancing units 106 when their consolidated health and load  
11 caches 1208 are unsynchronized with health and load tables 1204. This lack of  
12 synchronization may arise, for example, from a crash or other failure of  
13 consolidated health and load cache 1208 and/or of load balancing unit 106. As  
14 described above, each heartbeat message 1502 includes error checking data that is  
15 usable to verify equivalency between a consolidated health and load cache 1208  
16 and health and load tables 1204. If non-equivalency is discovered with regard to a  
17 particular host 108 and/or an application 316 thereof, the DNS name of the  
18 particular host 108 is acquired from the heartbeat messages 1502.

19 The DNS name is used by consolidated health and load cache 1208 to send  
20 a get table snapshot message 1508 to the particular host 108 in order to get  
21 updated health and load information 1206 in the form of a send table snapshot  
22 message 1510. A different or the same get table snapshot message 1508 is sent to  
23 each host 108 for which non-equivalency is discovered. Eventually, the health and  
24 load information 1206 in the consolidated health and load cache 1208 is equivalent  
25 to the health and load information 1206 in health and load tables 1204 as verifiable

1 by new heartbeat messages 1502. In this manner, a failed consolidated health and  
2 load cache 1208 can be bootstrapped back into operation without manual oversight  
3 using message protocol 1500 and an equivalency-checking scheme.

4 FIG. 17A and FIG. 17B illustrate exemplary health and load information  
5 proxy storage scenarios for health and load tables 1204 and for consolidated health  
6 and load caches 1208, respectively. In implementations described above with  
7 reference to FIGS. 12-16, hosts 108 include health and load infrastructure 1202.  
8 However, other implementations may entail hosts that do not include health and  
9 load infrastructure 1202.

10 For example, a host may be running a version of application(s) and/or an  
11 operating system for which health and load infrastructure is either not  
12 implemented or for policy reasons may not be installed on the host. Consequently,  
13 these types of hosts do not have health and load infrastructure 1202 executing  
14 thereon. Host 1702 is such a host that does not execute health and load  
15 infrastructure 1202. Nevertheless, host 1702 can utilize a health and load  
16 infrastructure 1202 that is executing on one or more proxies, such as proxy 1704.

17 Proxy 1704 has resident thereat and executing thereon a health and load  
18 infrastructure 1202, which includes a health and load table 1204. Host 1702 can  
19 use the functionality of health and load infrastructure 1202 by providing health  
20 and load information 1206 to health and load table 1204 for applications that are  
21 running on host 1702. Alternatively, proxy 1704 can deduce health and load on  
22 host 1702 by performing external monitoring actions. Proxy 1704 is illustrated as  
23 proxy 1704(1) and 1704(2) for redundancy and the resulting high availability.

24 In implementations described above with reference to FIGS. 12-16 and  
25 below with reference to FIG. 18, load balancing is effectuated with load balancing

1 units 106 that include consolidated health and load caches 1208. However, other  
2 implementations may entail load balancing that does not include consolidated  
3 health and load caches 1208.

4 For example, load balancing may be effectuated by monolithic load  
5 balancing hardware or other load balancing infrastructure that does not and/or  
6 cannot store or otherwise include a consolidated health and load cache 1208. Load  
7 balancer 1706 reflects such a load balancing device or devices that do not have a  
8 consolidated health and load cache 1208. Nevertheless, load balancer 1706 can  
9 utilize a consolidated health and load cache 1208 that exists on one or more  
10 proxies, such as proxy 1708.

11 Proxy 1708 includes a consolidated health and load cache 1208, which  
12 stores health and load information 1206 for hosted applications being serviced by  
13 load balancer 1706. Load balancer 1706 can use the health and load information  
14 1206 of consolidated health and load cache 1208 when performing load balancing  
15 functions by accessing such information using application programming interfaces  
16 (APIs) native to and supported by load balancer 1706. Alternatively, consolidated  
17 health and load cache 1208 can invoke APIs to push health and load information  
18 1206, including directives, to load balancer 1706. Proxy 1708 is illustrated as  
19 proxy 1708(1) and 1708(2) for redundancy and the resulting high availability.

20 FIG. 18 illustrates an exemplary target application endpoint allotment  
21 procedure that involves a classifier 304 and a health and load handler 314 of a load  
22 balancing unit 106. After health and load handler 314 has acquired a consolidated  
23 health and load cache 1208, health and load information 1206 thereof is utilized in  
24 the selection of application endpoints for new requests/connections.  
25

1 As described above with reference to FIG. 13B, consolidated health and  
2 load cache 1208 includes cached health and load information 1206 for multiple  
3 hosts 108. To facilitate the creation and updating of consolidated health and load  
4 cache 1208 from health and load information 1206 that originates from multiple  
5 hosts 108, the health and load information 1206 therein is organized so that it may  
6 be accessed by identifier of each host 108. However, the health and load  
7 information 1206 therein is also organized such that it can be accessed by type of  
8 application 316 in order to facilitate application endpoint selection.

9 In other words, health and load handler 314 is capable of accessing health  
10 and load information 1206 on a per-application 316 basis across health and load  
11 information 1206 for multiple hosts 108. Once health and load information 1206  
12 for a given application 316 has been accessed for each host 108, allocation of  
13 incoming connection requests may be performed in accordance with such health  
14 and load information 1206. For example, possible endpoints for the given  
15 application 316 may be allocated to incoming connection requests by selection of  
16 the endpoints of the given application 316 with consideration of available relative  
17 load capacity among healthy endpoints for the given application 316.

18 In a described implementation, classifier 304 makes a target application  
19 endpoint allotment request 1802 to health and load handler 314. As illustrated,  
20 target application endpoint allotment request 1802 includes (i) a virtual IP address  
21 and port, (ii) a protocol, and (iii) protocol-specification information. Target  
22 application endpoint allotment request 1802 therefore identifies a type of  
23 application 316 to which incoming connection requests are directed.

24 Health and load handler 314 receives target application endpoint allotment  
25 request 1802 and selects at least one physical endpoint corresponding to the

1 identified type of application 316 using any one or more of many selection  
2 mechanisms. To reduce latency, health and load handler 314 selects an allotment  
3 of application endpoints to be used over a number of incoming connection  
4 requests. This allotment is provided from health and load handler 314 to classifier  
5 304 using target application endpoint allotment response 1804. As illustrated,  
6 target application endpoint allotment response 1804 includes an allotment of  
7 physical IP addresses and ports (such as endpoints IP1, IP2, and IP3) for the  
8 identified type of application 316.

9 The allotment for target application endpoint allotment response 1804 may  
10 be completed using one or more allotment schemes. By way of example, a token  
11 allotment scheme 1806 and a percentage allotment scheme 1808 are illustrated.  
12 Token allotment scheme 1806 is a unit-based allotment scheme, and percentage  
13 allotment scheme 1808 is a time-based allotment scheme.

14 Token allotment scheme 1806 allocates tokens for each healthy endpoint  
15 IP1, IP2, and IP3 responsive to their relative load and capacity ratios. For the  
16 example as illustrated, of the total available capacity, IP1 has 40% of the available  
17 capacity, IP2 has 35% of the available capacity, and IP3 has 25% of the available  
18 capacity. Thus, the total number of tokens is divided along these percentages. The  
19 total number of tokens may be provided as part of target application endpoint  
20 allotment request 1802 or determined by health and load handler 314.

21 Any value for the total number of tokens may be used, such as 10, 45, 100,  
22 250, 637, 1000, and so forth. This value may be set in dependence on the number  
23 of connection requests per second and the speed/frequency at which application  
24 health and/or load is changing. Classifier 304 “uses up”/consumes one token  
25 when responding to each connection request with an application endpoint

1 allocation until the tokens are exhausted; classifier 304 then requests another token  
2 allotment using target application endpoint allotment request 1802.

3 Percentage allotment scheme 1808 determines available relative capacity in  
4 a similar manner. However, instead of tokens, these determined available relative  
5 capacities per application endpoint are provided to classifier 304 along with a  
6 duration timer 1810. Classifier 304 allocates target application endpoints to  
7 incoming connection requests in accordance with these available relative capacity  
8 percentages until expiration of duration timer 1810.

9 For percentage allotment scheme 1808, classifier 304 maintains a running  
10 record of application endpoint allocations to adhere to the allotted percentages and  
11 keeps track of time for duration timer 1810. When the timer expires, classifier 304  
12 then requests another percentage allotment using target application endpoint  
13 allotment request 1802.

14 It should be noted that token allotment scheme 1806 can also use a time  
15 limit. If allotted tokens are too old, they should be discarded and new ones  
16 acquired. Otherwise, classifier 304 may consume stale tokens that were  
17 previously allocated based on health and load information that is currently too  
18 outdated. Use of application endpoint allotments by classifier 304 is described  
19 further below in the section entitled "Exemplary Classifying, Forwarding, and  
20 Request Routing".

### 21 **Exemplary Session Tracking**

22 This section describes how host status information, such as session  
23 information, may be collected for and utilized in network load balancing. This  
24 section primarily references FIGS. 19-24 and illuminates session affinity  
25 preservation functionality such as that provided by session tracker 308 (of FIG. 3).

1 As described above with reference to FIGS. 1-3, each host 108 hosts one or more  
2 applications 316 that provide service(s) to clients 102. Session tracker 308 utilizes  
3 session information that relates to contexts for the connections established  
4 between applications 316 and clients 102 for certain described implementations of  
5 network load balancing.

6 FIG. 19 illustrates an exemplary network load balancing approach that  
7 involves session information 1902. At connection [1], client 102(1) is shown  
8 making a new connection with host 108(2) via load balancing infrastructure 106.  
9 Load balancing infrastructure 106 may be comprised of one or more load  
10 balancing units 106. When the connection request arrives at load balancing  
11 infrastructure 106, the request is typically routed to a host 108 using network load  
12 balancing functionality responsive to health and/or load information of hosts 108  
13 and/or applications 316 (not explicitly shown in FIG. 19) thereof.

14 When connection [1] is made, a session is established between client 102(1)  
15 and the servicing application 316, which is on host 108(2) in this example. The  
16 session provides a context for the communication exchange between client 102(1)  
17 and host 108(2). The information for the session context is stored at host 108(2).  
18 When connection [1] is completed, the session context may not be used again. On  
19 the other hand, the session context may be useful again if client 102(1) attempts to  
20 initiate another connection with hosts 108 for the service provided by application  
21 316. If this other connection is not routed to the same host 108(2) that stores that  
22 session context, then client 102(1) has to establish a new session context, which  
23 can be time consuming, data/processing intensive, and/or frustrating to the human  
24 user of client 102(1). With health and/or load information-based network load  
25



1 balancing, there is no likelihood greater than random chance that the second  
2 connection will be routed to 108(2).

3       However, if load balancing infrastructure 106 has access to a mapping  
4 between session information and hosts 108, load balancing infrastructure 106 can  
5 route connection requests that relate to previously established sessions to the  
6 appropriate host 108. Some session information may be inferred from the contents  
7 of packets flowing through load balancing infrastructure 106. However, this  
8 approach is imprecise and haphazard for a number of reasons. First, session  
9 establishment and termination is merely inferred. Second, some sessions are not  
10 “officially” terminated with an appropriate indication that is included in a packet.  
11 For example, some sessions simply time out. Third, packets being transmitted  
12 from host 108(2) to client 102(1) may take a path that does not include load  
13 balancing infrastructure 106, which precludes any snooping of such packets by  
14 load balancing infrastructure 106 for session information.

15       As shown in FIG. 19, hosts 108 provide session information (SI) 1902 to  
16 load balancing infrastructure 106. Using session information 1902 from hosts  
17 108, a session affinity preserver 1904 can preserve the affinity between an  
18 established session and the host 108 on which the session was established.  
19 Session information 1902 includes a linkage between or a mapping from each  
20 session established between a client 102 and a particular host 108 to that particular  
21 host 108. This mapping is accessible to session affinity preserver 1904 as part of  
22 host-session information mapping 1906. More-specific examples of session  
23 information 1902 are provided below especially with reference to FIGS. 20, 22,  
24 23A, and 23B.

1 In certain described implementations for session tracking, the logical nature  
2 of clients 102 is pertinent. As noted above with reference to FIG. 1, a client 102  
3 may be a specific device and/or a specific user of a device. Consequently, session  
4 affinity for a user client 102 that is accessing hosts 108 from different devices can  
5 still be preserved. Session continuations using session information 1902 can  
6 therefore still be effectuated in proxy scenarios (e.g., those of some internet  
7 service providers (ISPs)).

8 Continuing with the connection [1] example, the session established at host  
9 108(2) is provided to load balancing infrastructure 106 as session information  
10 1902. Specifically, a linkage/mapping between (i) the session context of client  
11 102(1) and host 108(2) and (ii) an identifier for host 108(2) is created at host-  
12 session information mapping 1906. When a connection request for connection [2]  
13 subsequently arrives for the same session context, session affinity preserver 1904  
14 locates this session context in host-session information mapping 1906 and  
15 ascertains that host 108(2) is associated with this session context from the  
16 linkage/mapping.

17 Responsive to the mapping of host 108(2) to the requested session context  
18 as ascertained by session affinity preserver 1904 from host-session information  
19 mapping 1906, connection [2] is routed to host 108(2). In this sense, preserving  
20 session affinity is a higher priority for load balancing infrastructure 106 than  
21 application health and load-based network load balancing decisions. However,  
22 health and/or load may be a more important network load balancing factor than  
23 session tracking when, for example, loading is extremely heavy or when the  
24 session-relevant application and/or host is in a failed condition.  
25

1 Many types of connections may be session-related. Examples include: a  
2 TCP connection, a transport layer security (TLS)/SSL session, a PPTP session, an  
3 IPSec/L2TP session, an ISA session, an HTTP cookie-based session, a Terminal  
4 Server session, an administrator-defined session, and so forth. By way of  
5 clarification, a TCP connection is considered to be a session of TCP packets. Also,  
6 a model for defining sessions by an administrator may be enumerated and  
7 supported. Furthermore, client IP-address-based sessions that are delineated by  
8 timeouts may also be supported. This is relatively non-intelligent session support,  
9 but is expected by some users.

10 A connection request from a client 102 varies by the type of desired  
11 session. For example, for sessions of type "TCP connection", the connection  
12 request comprises a TCP packet. For sessions of type "SSL session", the  
13 connection request comprises a TCP connection. Other such connection requests  
14 correspond to other session types. These examples also show how there may be  
15 session layers. At a lower session level, a session context for a TCP connection  
16 may include a TCP 4-tuple, a session number, the number of bytes sent/received,  
17 and so forth. At a higher session level, a session context for an SSL session may  
18 include a 32-byte session ID, a public key of the client 102 that is provided to the  
19 host 108, and so forth.

20 FIG. 20 illustrates an exemplary network load balancing approach that  
21 involves communicating session information using notifications 2006 and  
22 messages 2008. Multiple load balancing units 106(1), 106(2) ... 106(u) and  
23 multiple hosts 108(1), 108(2) ... 108(n) are shown. Each respective host 108(1),  
24 108(2) ... 108(n) includes one or more respective applications 316(1), 316(2) ...  
25 316(n) which are resident thereat and executing thereon. Notifications 2006 are

1 used to provide session information from applications 316, and messages 2008 are  
2 used to provide session information from hosts 108 to load balancing units 106.

3 As illustrated, each respective host 108(1), 108(2) ... 108(n) includes  
4 respective session tracking infrastructure (STI) 2002(1), 2002(2) ... 2002(n).  
5 Each respective session tracking infrastructure 2002(1), 2002(2) ... 2002(n)  
6 includes a respective session table 2014(1), 2014(2) ... 2014(n) (although only  
7 session table 2014(1) is explicitly illustrated in FIG. 19).

8 Each respective load balancing unit 106(1), 106(2) ... 106(u) includes  
9 respective traffic routing functionality (TRF) 2012(1), 2012(2) ... 2012(u). Traffic  
10 routing functionality 2012 may comprise, for example, classifying and/or  
11 requesting routing functionality, such as that provided by classifier 304 and  
12 request router 306, respectively. Distributed across load balancing units 106(1),  
13 106(2) ... 106(u) is a distributed session tracking manager 2010.

14 In a described implementation, traffic routing functionality 2012 and  
15 distributed session tracking manager 2010 are part of load balancing infrastructure  
16 106. Session tracking infrastructure 2002 may also be (e.g., a remote) part of load  
17 balancing infrastructure 106.

18 An API 2004 is employed to provide session information from applications  
19 316 to session tracking infrastructure 2002. Using API 2004, applications 316 are  
20 empowered to notify session tracking infrastructure 2002 of session information,  
21 including various changes thereto. More specifically, each application 316 is  
22 capable of providing, and session tracking infrastructure 2002 is capable of  
23 accepting, notifications 2006.

24 A notification that a session has been established (or session establishment  
25 notification 2006(E)) is provided from application 316 when a session is newly

1 established or opened. Session establishment notification 2006(E) includes a  
2 session identifier and optionally an identifier of application 316. A notification  
3 that a session has been terminated (or session termination notification 2006(T)) is  
4 provided from application 316 when a session is terminated or closed. Session  
5 termination notification 2006(T) also includes the session identifier and optionally  
6 the identifier of application 316.

7 When session tracking infrastructure 2002 accepts a session establishment  
8 notification 2006(E), it inserts an entry in session table 2014 for the new session.  
9 An exemplary session table 2014 is described further below with reference to FIG.  
10 23A. When session tracking infrastructure 2002 accepts a session termination  
11 notification 2006(T), it removes the entry in session table 2014 for the old session.

12 Session table 2014(1) is the authoritative source for session information  
13 1902 with respect to applications 316(1) on host 108(1). There is generally too  
14 much latency, however, to require traffic routing functionality 2012 to contact  
15 hosts 108 for access to session tables 2014 upon receipt of each incoming  
16 connection request having a session reference. Session information 1902 is  
17 therefore cached at load balancing units 106.

18 At load balancing units 106, distributed session tracking manager 2010  
19 caches session information 1902 as part of its session tracking management  
20 responsibilities. Generally, distributed session tracking manager 2010 is a  
21 distributed application and/or virtual service that resides partially on each load  
22 balancing unit 106. For each logical session, distributed session tracking manager  
23 2010 keeps at least one cached copy of session information therefor in a reliable  
24 and scalable manner that may be quickly utilized for routing traffic as incoming  
25

1 connection requests that have a session reference are received by load balancing  
2 infrastructure 106.

3 Communications between hosts 108 and load balancing units 106 are  
4 effectuated with a reliable protocol that ensures that messages 2008 sent from a  
5 host 108 arrive at the intended load balancing unit 106. Each host 108 is bound to  
6 at least one specific load balancing unit 106 that is the intended load balancing  
7 unit 106 for messages 2008. This binding is created by assigning an IP address of  
8 a specific load balancing unit 106 to each host 108 for sending session-tracking  
9 messages 2008 between session tracking infrastructure 2002 and distributed  
10 session tracking manager 2010. To facilitate high availability of load balancing  
11 infrastructure 106, if a load balancing unit 106 fails, another load balancing unit  
12 106 assumes the IP address of the failed load balancing unit 106. Failure detection  
13 for IP address assumption may be accomplished using a heartbeat or another  
14 aliveness monitoring scheme.

15 Thus, messages 2008 communicate session information 1902 from session  
16 tracking infrastructure 2002 to distributed session tracking manager 2010. For  
17 example, when session tracking infrastructure 2002 accepts a session  
18 establishment notification 2006(E), it also sends a session up message 2008(U) to  
19 distributed session tracking manager 2010. Session up message 2008(U) includes  
20 the session identifier, a host identifier, and optionally other information. Contents  
21 for a session up message 2008(U) are described further below with reference to  
22 FIG. 23B with respect to information that may be stored for each session by an  
23 implementation of distributed session tracking manager 2010. When session  
24 tracking infrastructure 2002 accepts a session termination notification 2006(T), it  
25 also sends a session down message 2008(D) to distributed session tracking

1 manager 2010. Messages 2008 can be sent before, during, or after session  
2 tracking infrastructure 2002 appropriately modifies session table 2014 in response  
3 to notifications 2006.

4 FIG. 21 is a flow diagram 2100 that illustrates an exemplary method for  
5 network load balancing that involves communicating session information using  
6 notifications and messages. Flow diagram 2100 includes fifteen blocks 2102-  
7 2130. Although the actions of flow diagram 2100 may be performed in other  
8 environments and with a variety of software schemes, FIGS. 1-3 and 19-20 are  
9 used in particular to illustrate certain aspects and examples of the method.

10 For example, the actions of four blocks 2102-2104 and 2118-2120 are  
11 performed by an application 316, the actions of six blocks 2106-2110 and 2122-  
12 2126 are performed by session tracking infrastructure 2002, and the actions of five  
13 blocks 2112-2116 and 2128-2130 are performed by a distributed session tracking  
14 manager 2010. The actions of eight of these blocks 2102-2116 are primarily  
15 directed to opening a session, and the actions of seven of these blocks 2118-2130  
16 are primarily directed to closing a session.

17 At block 2102, a session is opened. For example, application 316 may  
18 open a session with a client 102. At block 2104, a session establishment  
19 notification is provided. For example, application 316 may provide a session  
20 establishment notification 2006(E) to session tracking infrastructure 2002 using  
21 API 2004 as a consequence of and/or in conjunction with opening the session.

22 At block 2106, the session establishment notification is accepted. For  
23 example, session tracking infrastructure 2002 may accept session establishment  
24 notification 2006(E) from application 316 in accordance with API 2004. At block  
25 2108, an entry in a session table is inserted. For example, session tracking

1 infrastructure 2002 may insert an entry in session table 2014 for the opened  
2 session. Examples of such insertion are described further below especially with  
3 reference to FIG. 23A. At block 2110, a session up message is sent. For example,  
4 session tracking infrastructure 2002 may send a session up message 2008(U) to  
5 distributed session tracking manager 2010 using a reliable communication  
6 protocol.

7 At block 2112, the session up message is received. For example,  
8 distributed session tracking manager 2010 may receive session up message  
9 2008(U) from session tracking infrastructure 2002 in accordance with the reliable  
10 communication protocol. At block 2114, a session information entry is created.  
11 For example, distributed session tracking manager 2010 may create a session  
12 information entry for cached session information 1902 at one or more load  
13 balancing units 106. Examples of such creating and subsequent adding are  
14 described further below especially with reference to FIGS. 22 and 23B.

15 At block 2116, network traffic is routed with the session information. For  
16 example, traffic routing functionality 2012 in conjunction with distributed session  
17 tracking manager 2010 may use cached session information 1902, including the  
18 created session information entry, to route incoming connection requests that have  
19 a session reference. An example of such traffic routing is described further below  
20 especially with reference to FIG. 24. Additional examples are described below in  
21 the section entitled "Exemplary Classifying, Forwarding, and Request Routing".

22 At block 2118, the session is closed. For example, application 316 may  
23 close the session with client 102. At block 2120, a session termination notification  
24 is provided. For example, application 316 may provide a session termination  
25



1 notification 2006(T) to session tracking infrastructure 2002 using API 2004 as a  
2 consequence of and/or in conjunction with closing the session.

3 At block 2122, the session termination notification is accepted. For  
4 example, session tracking infrastructure 2002 may accept session termination  
5 notification 2006(T) from application 316 in accordance with API 2004. At block  
6 2124, the entry in the session table is removed. For example, session tracking  
7 infrastructure 2002 may remove the entry in session table 2014 for the closed  
8 session. At block 2126, a session down message is sent. For example, session  
9 tracking infrastructure 2002 may send a session down message 2008(D) to  
10 distributed session tracking manager 2010 using the reliable communication  
11 protocol.

12 At block 2128, the session down message is received. For example,  
13 distributed session tracking manager 2010 may receive session down message  
14 2008(D) from session tracking infrastructure 2002 in accordance with the reliable  
15 communication protocol. At block 2130, the session information entry is  
16 destroyed. For example, distributed session tracking manager 2010 may destroy  
17 the session information entry at the cached session information 1902 at any load  
18 balancing units 106 that have the session information entry. Examples of such  
19 destroying and subsequent deleting are described further below especially with  
20 reference to FIGS. 22 and 23B.

21 FIG. 22 illustrates an exemplary approach to managing session information  
22 at multiple load balancing units 106. Each respective load balancing unit 106(1),  
23 106(2) ... 106(u) includes a respective part 2202(1), 2202(2) ... 2202(u) of a  
24 distributed atom manager (DAM) 2202. DAM 2202 is an exemplary  
25 implementation of distributed session tracking manager 2010. Each respective

1 DAM portion 2202(1), 2202(2) ... 2202(u) includes a respective part 2206(1),  
2 2206(2) ... 2206(u) of a DAM table (DAMT) 2206.

3 DAM 2202 is a distributed application or virtual service that manages  
4 session information 1902 in a reliable and scalable manner so that traffic routing  
5 functionality 2012 can use it to preserve session affinity. For example, traffic  
6 routing functionality 2012 can access DAM 2202 using an API (not specifically  
7 shown) to search or have searched DAMT 2206. Function calls 2204, operation of  
8 DAM 2202, and other aspects of FIG. 22 are described further below after the  
9 description of FIGS. 23A and 23B.

10 FIG. 23A is an exemplary session table 2014 as illustrated in FIG. 20.  
11 Session table 2014 includes “v” entries 2302(1), 2302(2) ... 2302(v). Each entry  
12 2302 is inserted by session tracking infrastructure 2002 responsive to a session  
13 establishment notification 2006(E) that is accepted from an application 316. Each  
14 entry 2302 is removed by session tracking infrastructure 2002 responsive to a  
15 session termination notification 2006(T) that is accepted from application 316.

16 As described above, each session establishment notification 2006(E)  
17 includes a session identifier and optionally an identifier of application 316. Each  
18 respective entry 2302(1), 2302(2) ... 2302(v) in session table 2014 includes  
19 respective fields of (i) session identifier 2302(1I), 2302(2I) ... 2302(vI) and (ii)  
20 session type and/or application 2302(1T), 2302(2T) ... 2302(vT).

21 Session type and/or application 2302(T) may be “TCP”, “IPSEC”,  
22 “Terminal Server,” “HTTP-cookie”, an application type as noted above, and so  
23 forth. Session identifier 2302(I) may be “<source IP address, source TCP port,  
24 destination IP address, destination TCP port >”, “Client IP = 172.30.189.122”,  
25 “User = ‘joe\_user’”, “Cookie = ‘{b7595cc9-e68b-4eb0-9bf1-bb717b31d447}’”,

1 another e.g. application-specific identification for a session, and so forth. For TCP  
2 connection/session types, session identifier 2302(I) may alternatively be a  
3 canonical version of the TCP 4-tuple (for IPv4 or IPv6). Other values for the  
4 fields of session identifier 2302(I) and application/session type 2302(T) may  
5 alternatively be used.

6 FIG. 23B is an exemplary distributed atom manager (DAM) table (DAMT)  
7 2206 as illustrated in FIG. 22. DAM table 2206 includes “w” entries 2304(1),  
8 2304(2) ... 2304(w). Each session information entry 2304 is created by DAM  
9 2202 responsive to a session up message 2008(U) that is received from session  
10 tracking infrastructure 2002. Each session information entry 2304 is destroyed  
11 responsive to a session down message 2008(D) that is received from session  
12 tracking infrastructure 2002. As described further below, session information  
13 entries 2304 of DAM tables 2206 may actually be manipulated by DAM 2202  
14 using function calls 2204.

15 As described above, session up message 2008(U) includes the session  
16 identifier, a host identifier, and optionally other information. Each respective  
17 session information entry 2304(1), 2304(2) ... 2304(w) in DAM table 2206  
18 includes respective fields of (i) key 2304(1K), 2304(2K) ... 2304(wK), (ii) data  
19 2304(1D), 2304(2D) ... 2304(wD), and (iii) metadata 2304(1M), 2304(2M) ...  
20 2304(wM). For example, values for key 2304(K) fields may be alphanumeric  
21 strings, and values for data 2304(D) fields may be binary bits. Values for key  
22 2304(K) may be binary bits, too.

23 Key 2304(K) may correspond to the session identifier 2302(I). Data  
24 2304(D) may correspond to the host identifier, such as a network address of the  
25 host 108 on which the session context exists. Metadata 2304(M) may correspond

1 to other, optional information. Examples of such metadata 2304(M) include data  
2 that is used internally by DAM 2202 to resolve atom collisions and to track atom  
3 aliveness (e.g., via a time-out mechanism). (This characterization of entries 2304  
4 as being atomic is described more fully in the following paragraph.) More  
5 specifically, metadata 2304(M) includes, among other things, the identity of the  
6 entity (e.g., the instance of traffic routing functionality 2012) that added the  
7 session information entry 2304 to the DAM table 2206.

8 In a described implementation, each session information entry 2304 is  
9 atomic in the sense that DAM 2202 may add, delete, copy, etc. the entries 2304 as  
10 a whole, but DAM 2202 does not ordinarily modify a portion of any whole entry  
11 2304. Thus, atomic entries 2304 are added, deleted, copied, otherwise  
12 manipulated, etc. across DAM tables 2206 by DAM 2202 in order to implement  
13 availability and scalability for a session affinity preservation implementation.

14 Function calls 2204 (of FIG. 22) are usable by DAM 2202 to manipulate the  
15 atomic entries 2304 of DAM table 2206. Function calls 2204 may be  
16 communicated from one load balancing unit 106 to one or more other load  
17 balancing units 106 in a point-to-point or a multicast manner. These function calls  
18 include add atom 2204(A), delete atom 2204(D), query atom 2204(Q), and return  
19 atom 2204(R).

20 Add atom 2204(A) takes the form AddAtom(key, data) and is used to add  
21 an atomic entry 2304 to one or more DAM tables 2206. Hence, an add atom  
22 2204(A) function call may be formulated as AddAtom(<session identifier>, host  
23 IP address). Delete atom 2204(D) takes the form DeleteAtom(key) and is used to  
24 delete an atomic entry 2304 at one or more DAM tables 2206. Delete atom  
25 2204(D) function calls may be directed at those DAM tables 2206 known to have

1 a copy of the session that is identified by the key 2304(K) or may be multicast to  
2 all DAM tables 2206 to ensure that any copies are deleted.

3 Query atom 2204(Q) takes the form QueryAtom(key) and is used by a  
4 particular DAM portion 2202 when a session identifier as referenced by an  
5 incoming connection request is not located in the particular local DAM table 2206  
6 of the particular DAM portion 2202. Query atom 2204(Q) function calls are sent  
7 to one or more (including possibly all) other DAM portions 2202. In response,  
8 each other DAM portion 2202 checks its local DAM table 2206 for the  
9 key/session identifier. If the key is located by another DAM portion 2202, this  
10 other DAM portion 2202 replies with a return atom 2204(R).

11 Return atom 2204(R) takes the form ReturnAtom(key, data) and is used to  
12 reply to a query atom 2204(Q) function call. Return atom 2204(R) function calls  
13 are used when a DAM portion 2202 has a requested atomic entry 2304 in its local  
14 DAM table 2206 as identified by a key 2304(K) specified in the query atom  
15 2204(Q) function call. Return atom 2204(R) function calls may be directed back  
16 to the DAM portion 2202 that issued the query atom 2204(Q) function call.

17 Add atom 2204(A) function calls are used in response to session up  
18 messages 2008(U) and/or to replicate an atomic entry 2304 to one or more other  
19 DAM tables 2206. Such replication may be for redundancy and/or scalability.

20 Delete atom 2204(D) function calls are used in response to session down  
21 messages 2008(D) and may also be sent to one or more other DAM tables 2206.  
22 After an atomic entry 2304 is deleted, the atomic entry 2304 may enter a “zombie”  
23 state such that it remains with DAM 2202, and optionally so that it is actually still  
24 stored with DAM table 2206 with a zombie indication in the metadata 2304(M)  
25 field of the atomic entry 2304.

1        Thus, once an atomic entry 2304 is deleted, it may stay on in DAM 2202  
2 and DAM table 2206 in a zombie state so that packets for this (now dead and  
3 closed) session are directed to the host 108 of the session context for proper,  
4 protocol-specific treatment. For example, TCP packets received after a TCP  
5 connection has been torn down are directed to the host 108 that terminated the  
6 connection. This host 108 can respond appropriately – perhaps by sending an RST  
7 or by resending a FIN-ACK. The time the atomic entry 2304 spends in this  
8 zombie state matches (as closely as reasonably possible) the protocol-specific dead  
9 time of the reliable communication protocol that is employed.

10        A query atom 2204(Q) function call is used to attain an atomic entry 2304  
11 when a first load balancing unit 106 receives an incoming connection request that  
12 references a session that is not stored in the local DAM table 2206 of the DAM  
13 2202 of the first load balancing unit 106. It should be noted that other DAM  
14 portions 2202 may be queried simultaneously in a broadcast query atom 2204(Q)  
15 function call or sequentially until a positive return atom 2204(R) function call is  
16 received.

17        A return atom 2204(R) function call is used by a DAM portion 2202 of a  
18 second load balancing unit 106 to provide an atomic entry 2304 to the DAM  
19 portion 2202 of the first load balancing unit 106, where the atomic entry 2304 has  
20 a key 2304(K) that is specified by the key/session identifier in a query atom  
21 2204(Q) function call, which was previously issued by the DAM portion 2202 of  
22 the first load balancing unit 106. It should be noted that other components, such  
23 as traffic routing functionality 2012, may also be capable of calling functions  
24 2204, especially a query atom 2204(Q) function call, in accordance with an API or  
25 similar.

1 DAM portions 2202 and DAM tables 2206 may be organized and managed  
2 in a myriad of manners. Exemplary manners relate to replication/redundancy,  
3 local caching upon acquisition, hashing for location selection, and so forth. Zero,  
4 one, two, or more levels of replication up to full replication may be employed.  
5 With a zero level of replication, each atomic entry 2304 is stored at the DAM 2202  
6 that receives a session up message 2008(U) therefor without replication to other  
7 DAM portions 2202.

8 With a first level of replication, each atomic entry 2304 is stored at the  
9 DAM 2202 that receives a session up message 2008(U) therefor, and it is also  
10 added (copied) to one other DAM portion 2202 using an add atom 2204(A)  
11 function call. This handles one level of failure for a load balancing unit 106.  
12 Similarly, with a second level of replication, each atomic entry 2304 is stored at  
13 the DAM 2202 that receives a session up message 2008(U) therefor, and it is also  
14 added to two other DAM portions 2202. Generally, the one, two, etc. other DAM  
15 portions 2202 to which a given DAM portion 2202 copies atomic entries 2304 is  
16 predetermined or selected at random. Third, fourth, etc. levels of replication may  
17 also be employed.

18 Furthermore, full replication may be employed by having each atomic entry  
19 2304 that is stored at the DAM 2202 that receives a session up message 2008(U)  
20 therefor also being added to every other DAM portion 2202. Several factors are  
21 impacted by selection of the replication level: As the replication level increases,  
22 availability increases and latency decreases. On the other hand, network traffic  
23 and memory usage both increase as the replication level increases.

24 When full replication is not employed, local caching upon acquisition may  
25 be. For example, when a DAM portion 2202 does not locate a referenced session

1 identifier in its part of DAM table 2206, the DAM portion 2202 issues a query  
2 atom 2204(Q) function call to attain the atomic entry 2304 associated with the  
3 referenced session identifier via a return atom 2204(R) function call. Instead of  
4 jettisoning the attained atomic entry 2304 after use thereof, the DAM portion 2202  
5 caches the attained atomic entry 2304 in its part of DAM table 2206. This option  
6 offers a tradeoff between the above-enumerated factors.

7 As another option when full replication is not employed, hashing for  
8 location selection may be. The first atomic entry 2304 for a session is stored at the  
9 DAM portion 2202 that receives the session up message 2008(U). Replicated  
10 copy or copies are sent via add atom 2204(A) function calls to specific DAM  
11 portion(s) 2202 using a hashing function. Of a total range of possible hash values,  
12 each DAM portion 2202 is assigned a subset thereof. Each session identifier is  
13 hashed using some hashing function to arrive at a hashing value. This hashing  
14 value is mapped to the assigned DAM portion(s) 2202. The DAM portion 2202  
15 that first added the atomic entry 2304 then replicates the atomic entry 2304 to the  
16 assigned DAM portion(s) 2202.

17 With hashing for location selection, at least one DAM portion 2202 that has  
18 a desired atomic entry 2304 locally cached at its DAM table 2206 is knowable  
19 from the session identifier. A query atom 2204(Q) function call can therefore be  
20 directed to the known DAM portion(s) 2202. This usually reduces network traffic  
21 and/or latency.

22 This hashing for location selection may be used with one, two, three, or  
23 more levels of replication with each range of hashing values mapping to one, two,  
24 three, etc. different DAM portions 2202, respectively. Additionally, hashing for  
25 location selection may be used with local caching upon acquisition.



1        FIG. 24 is a flow diagram 2400 that illustrates an exemplary method for  
2        managing session information at multiple load balancing units. Flow diagram  
3        2400 includes eight blocks 2402-2416. Although the actions of flow diagram  
4        2400 may be performed in other environments and with a variety of software  
5        schemes, FIGS. 1-3, 19, 20, 22, and 23B are used in particular to illustrate certain  
6        aspects and examples of the method.

7        At block 2402, an incoming connection request with a session reference is  
8        analyzed. For example, traffic routing functionality 2012 may receive an  
9        incoming connection request that references a previously-opened/established  
10       session of a particular type. At block 2404, a local DAM table is searched using  
11       the session reference. For example, for a given load balancing unit 106 and traffic  
12       routing functionality 2012, the DAM portion 2202 thereof may search its  
13       corresponding DAM table 2206 looking for the session reference.

14       At block 2406, it is determined if the session reference matches a key of the  
15       local DAM table. For example, DAM portion 2202 may search key fields  
16       2304(K) of multiple entries 2304 of DAM table 2206 to determine whether the  
17       session reference matches any values of the key fields 2304(K). If so, flow  
18       diagram 2400 continues at block 2412.

19       If, on the other hand, the session reference does not match any key, flow  
20       diagram 2400 continues at block 2408. At block 2408, a query atom function call  
21       is made. For example, DAM portion 2202 may make a query atom 2204(Q)  
22       function call that includes the session reference/identifier as the key. The query  
23       atom 2204(Q) function call may be sent to at least one other DAM portion 2202.  
24       The number, selection, order, etc. of possible destination DAM portions 2202 for  
25       query atom 2204(Q) may depend on the options (e.g., replication level, hashing for

1 location selection, local caching upon acquisition, point-to-point versus multicast,  
2 etc.) employed by DAM 2202.

3 At block 2410, a returned atom is received. For example, information from  
4 a returned atom 2204(R) function call that is issued by another DAM portion 2202  
5 may be received. The other DAM portion 2202 successfully located an atomic  
6 entry 2304 in its corresponding DAM table 2206, with the located atomic entry  
7 2304 having a key that matches the session reference. The information from the  
8 returned atom 2204(R) function call includes values from key field 2304(K) and  
9 data field 2304(D) for the located atomic entry 2304. These values correspond to  
10 the session identifier of the session and the network address of the host 108 that is  
11 affinitized to the session.

12 At block 2412, an atomic entry is extracted. The atomic entry is extracted  
13 from the local DAM table if a match was found locally (at blocks 2404 and 2406)  
14 or from the returned atom if a match was found elsewhere (at blocks 2408 and  
15 2410). For example, an atomic entry 2304 may be extracted from DAM table  
16 2206 of the DAM portion 2202 or from information received by a return atom  
17 2204(R) function call. The extracted atomic entry 2304 may be cached at the local  
18 DAM table 2206 if received as a result of the return atom 2204(R) function call.

19 At block 2414, the host having session affinity with the referenced session  
20 is ascertained from the atomic entry. For example, a value of the data field  
21 2304(D) of the extracted atomic entry 2304 may be ascertained to thereby  
22 ascertain a network address of the affinitized host 108. At block 2416, the  
23 incoming connection request is routed to the ascertained host. For example, traffic  
24 routing functionality 2012 and/or forwarding functionality may route the incoming  
25 connection request having the session reference to the ascertained and affinitized

1 host 108. Exemplary classifying, request routing, and forwarding functionalities  
2 are described in the following section.

### 3 **Exemplary Classifying, Forwarding, and Request Routing**

4 This section describes how traffic routing may be implemented for network  
5 load balancing, including with regard to high availability of such traffic routing  
6 functionality. Traffic routing functionality may include classifying and/or  
7 requesting routing functionality, especially in conjunction with forwarding  
8 functionality. This section primarily references FIGS. 25-31. It illuminates the  
9 functionality of a request router 306 (of FIG. 3), an interrelationship between  
10 tracking sessions and utilizing health and load information when routing traffic,  
11 operational implementations for traffic routing interactions with session  
12 information and/or health and load information, failover procedures for high  
13 availability of network load balancing infrastructure (including handling failures  
14 of classifying, forwarding, and/or request routing components), additional network  
15 load balancing infrastructure configurations, and so forth.

16 FIG. 25 illustrates exemplary network load balancing infrastructure having  
17 request routing functionality as realized by request router 306(H/S). As noted  
18 above with reference to traffic routing functionality 2012, traffic routing may  
19 relate to classifying (e.g., with forwarding) and/or requesting routing. Packet-level  
20 classifying, in conjunction with forwarding, is described above with particular  
21 reference to FIG. 4. Request routing is described here with particular reference to  
22 FIG. 25.

23 Request-level routing occurs at a higher level than that of packet-level  
24 routing. Generally, a request router 306 acts as a proxy for an application 316  
25 running on a host 108. Request router 306 terminates TCP connections, parses

1 (perhaps partially) each request from a client 102, and resubmits each request to  
2 host 108. Request router 306 may perform pre-processing on the connection, such  
3 as SSL decryption. Also, request router 306 may chose to absorb certain requests  
4 (e.g., the request router may maintain a cache of responses), and it may  
5 “arbitrarily” modify requests before forwarding them to hosts 108.

6 Request routers 306 are usually application-specific, and they may be rather  
7 open-ended in what they are capable of doing. By way of example only, a single  
8 class of request routers 306 – HTTP/SSL request routers 306(H/S) – are addressed  
9 in the following description. As illustrated, a client 102 having a network address  
10 C1 is communicating across network 104 with hosts 108(1) and 108(2) having  
11 network addresses H1 and H2, respectively. The communications are effectuated  
12 via load balancing infrastructure that includes an HTTP/SSL request router  
13 306(H/S).

14 HTTP/SSL request router 306(H/S) terminates HTTP and SSL traffic,  
15 decrypts SSL traffic, examines each HTTP request from client 102, applies  
16 application-specific rules to classify each request and to determine the “best”  
17 endpoint for that request while taking into account application endpoint health and  
18 load information, and submits the request to the endpoint. The request submission  
19 to the endpoint uses a separate TCP connection than that of the one originated by  
20 client 102 (the latter connection is terminated at HTTP/SSL request router  
21 306(H/S)). These actions may be considered as logically equivalent to the actions  
22 performed by a classifier 304, but a difference arises in that these actions in  
23 HTTP/SSL request router 306(H/S) are occurring at the logical request level for  
24 each request within the TCP connection. HTTP/SSL request router 306(H/S), and  
25

1 request routers 306 generally, can use the same (i) application health and load and  
2 (ii) session tracking infrastructure that is used by classifiers 304.

3 HTTP/SSL request router 306(H/S) is acting as an intermediary between  
4 client 102 and two hosts 108(1) and 108(2). It is handling two requests from client  
5 102 over a single TCP connection. In a described implementation, the resulting  
6 request routing involves a number of actions. First, client 102 establishes an http  
7 or https connection [1] to HTTP/SSL request router 306(H/S) and sends a request  
8 #1 2502(1).

9 Second, HTTP/SSL request router 306(H/S) terminates the SSL session (if  
10 the traffic is SSL encrypted), parses request #1 2502(1), and examines the content  
11 of request #1 2502(1). Taking into account application health and load as well as  
12 session information, HTTP/SSL request router 306(H/S) determines that host  
13 108(1) is the “best” host for this particular request #1 2502(1) in this example.

14 Third, HTTP/SSL request router 306(H/S) establishes a secondary TCP  
15 connection [2] to host 108(1). This secondary TCP connection is not sourced from  
16 a VIP address on network 104; instead, it is sourced from an address (not shown in  
17 FIG. 25) that is dedicated to request router 306(H/S) to ensure that responses 2504  
18 from host(s) 108 reach the correct request router 306. (There may be multiple  
19 request routers 306 that are active even though one request router 306(H/S) is  
20 shown in FIG. 25 for clarity.) It may alternatively use an existing connection [2] to  
21 host 108(1). HTTP/SSL request router 306(H/S) then sends an e.g. unencrypted  
22 version of request #1 2502(1) to host 108(1). Fourth, host 108(1) replies with a  
23 response #1 2504(1). Fifth, HTTP/SSL request router 306(H/S) encrypts this  
24 response #1 2504(1) and sends it back to client 102 on TCP connection [1].  
25

1 Sixth, client 102 sends another request, request #2 2502(2). Request #2  
2 2502(2) is handled similarly to the handling of request #1 2502(1), except that  
3 HTTP/SSL request router 306(H/S) selects host 108(2). The different selection  
4 may be because host 108(1) is now failing or more-heavily loaded, because  
5 request #2 2502(2) is directed to a different URL than request #1 2502(1), and so  
6 forth. Regardless, HTTP/SSL request router 306(H/S) establishes another  
7 secondary TCP connection, but this secondary TCP connection [3] is to host  
8 108(2). Unencrypted request #2 2502(2) is routed to host 108(2), and a response  
9 #2 2504(2) is received therefrom as a result. An encrypted version of response #2  
10 2504(2) is then sent to client 102 from HTTP/SSL request router 306(H/S).

11 Seventh, client 102 closes TCP connection [1] with HTTP/SSL request  
12 router 306(H/S). HTTP/SSL request router 306(H/S) (at some future time) closes  
13 connections [2] and [3] that it made to hosts 108(1) and 108(2), respectively, on  
14 behalf of client 102. TCP connection [2] may alternatively be closed after  
15 HTTP/SSL request router 306(H/S) decides to open/use TCP connection [3] for  
16 request #2 2502(2).

17 Because an HTTP/SSL request router 306(H/S) terminates the  
18 HTTP/HTTPS connection, HTTP/SSL request router 306(H/S) can do more than  
19 route requests. For example, HTTP/SSL request router 306(H/S) can potentially  
20 maintain its own cache of responses (e.g., with an out-of-band mechanism to  
21 invalidate the cache). As noted in the above example, HTTP/SSL request router  
22 306(H/S) can also potentially route different kinds of requests to different sets of  
23 hosts 108 based on e.g. the requested URL. Conversely, HTTP/SSL request router  
24 306(H/S) can potentially aggregate requests from many short-lived client  
25 connections and send them over a few, long-standing TCP connections to hosts

1 108. Such connection aggregation can reduce the TCP connection processing  
2 overhead in hosts 108.

3 Request routers of other classes may correspond to other exemplary  
4 protocols besides HTTP. For example, a request router may be a SOAP request  
5 router. SOAP request routers function analogously to an HTTP/SSL request router  
6 306(H/S). However, SOAP request routers specialize in routing SOAP traffic.  
7 SOAP request routers understand SOAP headers and make routing decisions based  
8 on the SOAP headers as well as application health and load.

9 Both packet-level classification and forwarding (or packet-level routing)  
10 and request-level routing can provide some form of layer-7 load balancing. Layer-  
11 7 load balancing is described further below in the section entitled "Exemplary  
12 Connection Migrating with Optional Tunneling and/or Application-Level Load  
13 Balancing". Packet-level routing provides read-only access to the initial portion of  
14 a client's TCP connection data, and request-level routing provides read and modify  
15 access to an entire data stream.

16 Packet-level routing typically has several advantages over request-level  
17 routing. These advantages include transparency (client packets are delivered to  
18 hosts as-is, preserving source and destination IP addresses and port numbers), low  
19 processing overhead (generally, forwarding traffic involves a route lookup), low  
20 latency (individual packets are forwarded, and packets are not queued once the  
21 TCP connection destination has been determined), and high-availability (generally,  
22 a failure in a forwarder does not terminate the TCP connection). Request-level  
23 routing, on the other hand, typically has the following advantages over packet-  
24 level routing: an ability to examine an entire data stream flowing to and from the  
25

1 client; and an ability to transform a data stream, and even to split the data stream  
2 among multiple hosts or aggregate data streams from multiple clients.

3 FIG. 26 is a flow diagram 2600 that illustrates an exemplary method for  
4 routing incoming packets with regard to (i) session information and (ii) health and  
5 load information. Flow diagram 2600 includes eight blocks 2602-2616. Although  
6 the actions of flow diagram 2600 may be performed in other environments and  
7 with a variety of software schemes, FIGS. 1-3, 12, 18-20, 22, and 23B are used in  
8 particular to illustrate certain aspects and examples of the method.

9 At block 2602, an incoming packet is received. For example, a packet from  
10 a client 102 may be received at a forwarder 302 of a load balancing unit 106. At  
11 block 2604, it is determined if the received packet is for a preexisting session. For  
12 example, forwarder 302 may consult a local DAM table 2206( ) to determine that  
13 the received packet is already part of a TCP/IP session.

14 Additionally, forwarder 302 may consult the local DAM table 2206( ) and  
15 determine that the received packet is not already part of a TCP/IP session. In this  
16 case, forwarder 302 provides the received packet to a classifier 304, which checks  
17 for a higher level session affinity for the received packet if it has a session  
18 reference. Examples for these actions are described above with particular  
19 reference to FIG. 24 and further below with particular reference to FIGS. 27 and  
20 28.

21 If the received packet is for a preexisting session (as determined at block  
22 2604), then flow continues at block 2606. At block 2606, a host that is affinitized  
23 to the preexisting session is ascertained. For example, an affinitized host 108 may  
24 be ascertained from the local DAM 2206( ) and/or the overall distributed DAM  
25 2206 by forwarder 302 or classifier 304.



1       At block 2608, it is determined if the affinitized host is healthy. For  
2       example, classifier 304 may consult a consolidated health and load cache 1208 to  
3       determine if the affinitized host 108 is healthy, especially for those received  
4       packets that are part of sessions that are of a higher logical level than TCP/IP  
5       sessions. The action(s) of this block may be accomplished in conjunction with a  
6       health and load handler 314.

7       If the affinitized host is healthy (as determined at block 2608), then flow  
8       continues at block 2610. At block 2610, the received packet is routed to the  
9       affinitized host. For example, forwarder 302 (for TCP/IP sessions) or classifier 304  
10      (for higher-level sessions) may route the packet to the affinitized host 108. In an  
11      alternative implementation, classifier 304 may return the received packet to  
12      forwarder 302 for routing to the affinitized host 108 even for received packets that  
13      are part of higher-level sessions.

14      If, on the other hand, the affinitized host is not healthy (as determined at  
15      block 2608), then flow continues at block 2612. Also, if on the other hand, the  
16      received packet is not for a preexisting session (as determined at block 2604), then  
17      flow continues at block 2612. At block 2612, a host is selected responsive to  
18      health and load information. For example, classifier 304 may select a host 108  
19      from and/or using a health and load-related application allotment (e.g., from a  
20      target application endpoint allotment response 1804) that is attained from health  
21      and load handler 314. Examples for these action(s) are described above with  
22      particular reference to FIGS. 19 and 18 and further below with particular reference  
23      to FIG. 30.

24      At block 2614, the received packet is routed to the selected host. For  
25      example, classifier 304 may route (optionally via forwarder 302) the packet to the

1 selected host 108. At block 2616, a route for a connection path to the selected host  
2 is plumbed. For example, classifier 304 may add a session information entry to  
3 DAM table 2206, especially at the DAM table 2206( ) that is local to the forwarder  
4 302 that provided the received packet to the classifier 304. This session  
5 information entry may be replicated in accordance with the instituted redundancy  
6 policy for a DAM 2202 (e.g., of a session tracker 308).

7 The action(s) of block 2614 and those of block 2616 may be performed in  
8 the order specifically illustrated, with those of block 2616 being performed prior  
9 to those of block 2614, with the actions partially or fully overlapping in any order,  
10 and so forth. It should be noted that the actions performed by classifier 304 as  
11 described above may alternatively be performed by a request router 306 (or more  
12 generally traffic routing functionality 2012).

13 In addition to packet-level and request-level routing, traffic routing  
14 functionality as described herein (e.g., traffic routing functionality 2012, a request  
15 router 306, a forwarder 302/classifier 304 pair, etc.) can also be used to implement  
16 firewall functionality. Hence, a feature of the traffic routing functionality may  
17 include blocking traffic, instead of automatically routing traffic to the correct host  
18 108. For example, a classifier 304 can inspect traffic and drop it if it is deemed  
19 unsafe.

20 FIG. 27 illustrates an exemplary traffic routing flow in the absence of  
21 failures. As illustrated, one or more load-balancing-aware switches 202(LBA)  
22 front the remaining load balancing infrastructure 106 (not separately indicated).  
23 Forwarding and classifying functionality are distributed across three devices or  
24 nodes. A first device includes forwarder 302(1) and classifier 304(1). A second  
25 device includes classifier 304(2). A third device includes forwarder 302(2).

1 With classifier 304(2) executing on the second device and forwarder 302(2)  
2 executing on the third device, each device may be specially tuned for its respective  
3 functions. For example, the hardware, software, firmware, some combination  
4 thereof, etc. of the second device and the third device may be adapted to support  
5 the desired functionality without excessive over provisioning. Thus, the third  
6 device that includes forwarder 302(2) may be akin to a switch and/or router from a  
7 hardware capability perspective, and the second device that includes classifier  
8 304(2) may be more akin to a server and/or personal computer from a hardware  
9 capability perspective.

10 Although shown as three devices that are providing functionality across  
11 four components, alternative logical and/or device-level configurations for  
12 forwarding and classifying functionality are applicable to the exemplary traffic  
13 routing flow that is described here for FIG. 27. Also, although the routing  
14 destinations are shown as hosts 108, the descriptions herein of routing  
15 implementations may alternatively be applied more generally to a next node  
16 destination for the packet and not necessarily a final node that consumes the  
17 packet.

18 A DAM 2202 realization of session tracker 308 is used to implement DAM  
19 table 2206. However, session affinity preservers 1904 in general are also  
20 applicable to the exemplary traffic routing flow of FIG. 27. Forwarder 302(1)  
21 includes DAM table portion 2206(1), and forwarder 302(2) includes DAM table  
22 portion 2206(2). Incoming packets are routed to host 108(1) or host 108(2).

23 In a described implementation, DAM 2202 is a distributed, in-memory  
24 table of “atoms” 2304 (e.g., keyword-value pairs, with optional metadata) having  
25 session information. DAM 2202 and DAM table 2206 is described further above

1 with particular reference to FIGS. 22-24. Any node in the cluster of classifiers 304  
2 may add, query, and delete atoms 2304. DAM 2202 maintains a highly available  
3 DAM table 2206 that includes active (e.g., TCP/IP level) routes as well as higher-  
4 level session information. Examples of higher level sessions include: a TLS/SSL  
5 session, a PPTP session, an IPSec/L2TP session, an ISA session, an HTTP cookie-  
6 based session, and so forth. Furthermore, DAM 2202 may include session  
7 information entries in DAM table 2206 that are directed to other non-TCP/IP  
8 sessions, such as RTP, UDP, and so forth.

9 At (1), load-balancing-aware switches 202(LBA) direct an incoming packet  
10 to forwarder 302(1). At (2), forwarder 302(1) consults its internal routing table,  
11 DAM table 2206(1). When forwarder 302(1) does not find an atomic entry 2304  
12 for this packet, it forwards the packet to its assigned and/or associated classifier,  
13 classifier 304(1).

14 At (3), classifier 304(1) recognizes that the packet in this example is a first  
15 packet of a new session (e.g., a SYN packet for a TCP connection). Classifier  
16 304(1) therefore treats the packet as a start of a new TCP connection from a client  
17 102. Using health and load information from a health and load handler 314 (not  
18 explicitly illustrated), classifier 304(1) determines that host 108(1) should receive  
19 this session.

20 Classifier 304(1) updates DAM table 2206(1) that serves as the local  
21 routing table for forwarder 302(1), and it also inserts an atomic entry 2304  
22 representing the route into the overall DAM 2206. These may be separate  
23 operations, a single operation in which the TCP/IP-level sessions of DAM table  
24 2206 are located at forwarders 302, and so forth. DAM 2202 internally replicates  
25 this route to one or more other members of the cluster of classifiers 304 in

1 accordance with its stipulated redundancy policy. Classifier 304(1) may  
2 optionally communicate with host 108(1) to confirm the creation of the new  
3 session before it updates DAM table 2206(1) of forwarder 302(1) and the overall  
4 DAM 2202/DAM table 2206.

5 At (4), forwarder 302(1) directly forwards subsequent packets for this  
6 connection to host 108(1) without interacting with classifier 304(1). DAM 2202  
7 can be used to mask, at least in part, the failure of a forwarder 302, a classifier  
8 304, or a forwarder/classifier pair 302/304. DAM 2202 can also be used, at least  
9 in part, to preserve client connectivity if load-balancing-aware switches 202(LBA)  
10 inadvertently start sending packets for an established connection to a different  
11 forwarder 302.

12 FIG. 28 illustrates an exemplary traffic routing flow in the presence of  
13 failure(s). In contrast to the “failure-free” exemplary traffic routing flow of FIG.  
14 27, a failure has occurred in a portion of the network load balancing infrastructure  
15 106 (not specifically identified) of FIG. 28. Specifically, the first device, on which  
16 forwarder 302(1) and classifier 304(1) are resident and executing, fails after the  
17 connection that is illustrated in FIG. 27 is established. This failure is masked, at  
18 least in part, by DAM 2202.

19 At (1), load-balancing-aware switches 202(LBA) detect the failure of  
20 forwarder 302(1) and start forwarding packets for the connection to some other  
21 forwarder 302 in the cluster. In this example, the other forwarder 302 is forwarder  
22 302(2). Although FIG. 28 illustrates a failure situation, load-balancing-aware  
23 switches 202(LBA) may also send this traffic to forwarder 302(2) even if  
24 forwarder 302(1) is still available. This non-failure-induced change of forwarders  
25 302 may occur, for example, because load-balancing-aware switches 202(LBA) do

1 not preserve the affinity of this traffic to forwarder 302(1). Any of several factors  
2 can cause switches 202 to (mis)direct traffic to a different, non-affinitized  
3 forwarder 302. For example, traffic for the same higher-level session can arrive at  
4 switches 202 from a different source IP address or source port when the source is  
5 behind a farm of proxy servers. The actions of notations (2)-(5) apply to both the  
6 failure and the “misdirected traffic” situations.

7 At (2), forwarder 302(2) consults its routing table, DAM table 2206(2).  
8 When it does not find a route for this packet, it forwards the packet to its classifier  
9 304(2). At (3), classifier 304(2) recognizes that this packet is a “mid-session”  
10 packet, and classifier 304(2) queries DAM 2202 for the route for this packet.  
11 DAM 2202 responds with the route for the connection from an atomic entry 2304  
12 that is associated therewith.

13 At (4), classifier 304(2) plumbs the route in forwarder 302(2). An  
14 exemplary protocol for plumbing routes is described further below. At (5),  
15 subsequent packets for this connection that are directed to forwarder 302(2) are  
16 routed directly to the correct host, which is host 108(1) in this example, without  
17 consulting classifier 304(2).

18 Generally, a route plumbing protocol for communications between  
19 classifiers 304 and forwarders 302 includes instructions to add and remove routes.  
20 More specifically, an add route instruction is sent from a classifier 304 to a  
21 forwarder 302 in order to plumb a route from the forwarder 302 to a destination  
22 host 108 for a given connection. By way of example, an add route instruction can  
23 be provided to forwarder 302(2) from classifier 304(2) as indicated at (4) in FIG.  
24 28. The route (e.g., a key and corresponding value) is added to local DAM table  
25 2206(2) for quick access by forwarder 302(2) in the future. In this example,

1 classifier 304(2) is a separate device from forwarder 302(2), so the route plumbing  
2 protocol may be an inter-device protocol. However, the route plumbing protocol  
3 may also be utilized for intra-device communications.

4 In a described implementation, classifier 304(2) includes a connection  
5 inventory 2802. With connection inventory 2802, classifier 304(2) keeps track of  
6 the sessions of any forwarders 302 (such as forwarder 302(2)) for which classifier  
7 304(2) plumbs routes. To enable classifier 304(2) to keep track of the sessions,  
8 including cessations thereof, forwarder 302(2) forwards final packets for sessions  
9 (such as a TCP FIN packet) to classifier 304(2). Classifier 304(2) then deletes an  
10 entry in connection inventory 2802 that corresponds to the session and sends a  
11 delete route instruction to forwarder 302(2). Upon receiving the delete route  
12 instruction, forwarder 302(2) removes the corresponding route in DAM table  
13 2206(2).

14 In this manner, the classifying functionality in conjunction with session  
15 tracking functionality can control the route tables, and the routes thereof, that are  
16 used by the forwarding functionality. Consequently, forwarding functionality that  
17 is separated onto a different device may be effectuated using high-speed, but  
18 relatively simple, hardware. Alternatively, classifiers 304 may rely on  
19 communications with/from hosts 108, rather than (or in addition to) intercepted  
20 session initiation (such as TCP SYN) and termination (such as TCP FIN) packets,  
21 to determine the lifetimes of sessions. In other words, classifiers 304 may  
22 alternatively or additionally receive and utilize session (up/down) messages  
23 2008(U/D) as described above in the section entitled "Exemplary Session  
24 Tracking".  
25

1        FIG. 29 illustrates additional exemplary failover procedures for high  
2        availability of network load balancing infrastructure 106. Failover procedures for  
3        two different failures, failure 2902 and failure 2906, are described. As illustrated,  
4        network load balancing infrastructure 106 (not separately indicated) includes five  
5        components: forwarder 302(1), forwarder 302(2), forwarder 302(3), classifier  
6        304(1), and classifier 304(2).

7        In a described implementation, each of these five components 302(1),  
8        302(2), 302(3), 304(1), and 304(2) corresponds to an individual device. However,  
9        similar failover procedures apply to environments in which different load  
10       balancing components share devices. Also, similar or analogous failover  
11       procedures may apply to environments having other numbers, combinations,  
12       scalings, etc. of components.

13       Initially at [1], router/switch(es) 202 direct an incoming packet that  
14       happens to be for a new connection to forwarder 302(1). Because forwarder  
15       302(1) does not have a route for this connection in its local routing table, it sends  
16       the packet to classifier 304(1) as indicated by the dashed double arrow at (1).  
17       Classifier 304(1) first checks session information with reference to session  
18       tracking 308 for a possible higher-level session affinity. In this example, the  
19       packet is not affinized to an existing session, so classifier 304(1) selects a host 108  
20       with reference to health and load information with reference to health and load  
21       handling 314.

22       Specifically, classifier 304(1) selects host 108(1) in this example.  
23       Assuming the packet is for a TCP/IP connection, this TCP/IP session as linked to  
24       host 108(1) is added to DAM 2202 using an add atom 2204(A) function call by  
25       classifier 304(1). The initial packet is forwarded to host 108(1) by classifier



1 304(1) or forwarder 302(1). Classifier 304(1) also plumbs a route in the local  
2 routing table of forwarder 302(1). Subsequent packets are forwarded to host  
3 108(1) by forwarder 302(1) without further interaction with classifier 304(1).

4 At some time during connection [1], there is a failure 2902 at forwarder  
5 302(1). With load-balancing-aware router/switch(es) 202(LBA), this failure 2902  
6 is detected. As a result, at point 2904, router/switch(es) 202 direct later packets  
7 that would have been sent to forwarder 302(1) along connection [1] to another  
8 forwarder 302, which is forwarder 302(2) in this example.

9 Forwarder 302(2) thus receives future packets along a connection [2].  
10 Because forwarder 302(2) does not have an entry in its local routing table for the  
11 packets that were formerly directed to forwarder 302(1), forwarder 302(2) sends  
12 the first received packet of connection [2] to the classifier to which it is  
13 assigned/associated. In this example, forwarder 302(2) is assigned to classifier  
14 304(2) as indicated by the dashed double arrow at (2).

15 Classifier 304(2) uses a query atom 2204(Q) function call to attain the  
16 atomic entry 2304 (not explicitly shown) from DAM 2202 that is associated with  
17 the existing TCP/IP connection. This atomic entry 2304 is provided through DAM  
18 2202 of session tracking 308 via a return atom 2204(R) function call. Classifier  
19 304(2) extracts the host 108(1) that is affinitized with this TCP/IP connection from  
20 the returned atomic entry 2304. Classifier 304(2) forwards the first received  
21 packet for connection [2] to host 108(1) and also plumbs a route in the local  
22 routing table of forwarder 302(2). Subsequent packets are forwarded to host  
23 108(1) by forwarder 302(2) without further interaction with classifier 304(2).

24 The above descriptions focus predominantly on failures of individual  
25 forwarder 302 components. However, classifier 304 components can also fail.

1 For example, at some point, there is a failure 2906 at classifier 304(2). Forwarder  
2 302(2) detects failure 2906 when it attempts to consume classification services or  
3 through noticing a lack of some aliveness indication such as a heartbeat-type  
4 indicator. To handle failure 2906, forwarder 302(2) is reassigned or re-associated  
5 with a different classifier 304, which is classifier 304(1) in this example. Future  
6 classification functionality is provided to forwarder 302(2) by classifier 304(1) as  
7 indicated by the dashed double arrow at (3).

8 FIG. 30 illustrates an exemplary operational implementation of traffic  
9 routing interaction with health and load information. Forwarder 302 and classifier  
10 304 interact with health and load handler 314 in order to route packets to hosts  
11 108(1), 108(2) ... 108(n). Although a forwarder 302 and a classifier 304 are  
12 illustrated, the exemplary operational implementation is also applicable to a  
13 request router 306 (or traffic routing functionality 2012 in general).

14 As illustrated, host 108(1) includes application endpoints IP1, IP3, and IP4  
15 for application #1, application #1, and application #2, respectively. Host 108(2)  
16 includes application endpoints IP2 and IP6 for application #1 and application #2,  
17 respectively. Host 108(n) includes application endpoint IP5 for application #2.  
18 These hosts 108(1), 108(2) ... 108(n) and application endpoints IP1, IP2, IP3, IP4,  
19 IP5, and IP6 are monitored by health and load handler 314 (e.g., using health and  
20 load infrastructure 1202, consolidated health and load cache 1208, etc.).

21 In a described implementation, at (1) classifier 304 requests one or more  
22 application endpoint allotments (e.g., via at least one target application endpoint  
23 allotment request 1802) in an environment using a token allotment scheme 1806.  
24 Health and load handler 314, in this example, responds by providing token  
25

1 allotments 3002 (e.g., via at least one target application endpoint allotment  
2 response 1804).

3 Specifically, a token allotment for application #1 3002(1) and a token  
4 allotment for application #2 3002(2) are available to classifier 304. Token  
5 allotment for application #1 3002(1) initially provides 40 tokens for IP1, 35 tokens  
6 for IP2, and 25 tokens for IP3. Token allotment for application #2 3002(2)  
7 provides 10 tokens for IP4, 72 tokens for IP5, and 18 tokens for IP6. For each  
8 new connection that is allocated a routing to an application endpoint by classifier  
9 304, a token is consumed by classifier 304.

10 At (2), forwarder 302 receives an initial incoming packet for a new  
11 connection. Because no routing for this new connection is present in local DAM  
12 table portion 2206 of forwarder 302, forwarder 302 forwards the initial packet to  
13 classifier 304 at (3).

14 At (4), classifier 304 (e.g., after determining that the initial packet does not  
15 include a session reference for a higher-level session) selects an application  
16 endpoint (and thus a host 108) responsive to health and load information.  
17 Specifically, for a new connection that is to be served by application #1, classifier  
18 304 can select any of IP1, IP2, and IP3 if a token for the respective endpoint still  
19 exists.

20 Classifier 304 can consume tokens in any of many possible manners. For  
21 example, classifier 304 may use a round-robin approach regardless of the number  
22 of tokens per endpoint. Alternatively, classifier 304 may simply start from IP1 and  
23 progress through IP3 while consuming all tokens for each endpoint before moving  
24 to the next endpoint in a linear approach. Also, classifier 304 may consume a  
25 token from the endpoint-defined-set of tokens that currently has the greatest

1 number of tokens at any one moment. Using the latter approach, classifier 304  
2 selects IP1. Other approaches may also be employed.

3 As illustrated, classifier 304 consumes a token for application endpoint IP2.  
4 Consequently, the token set for IP2 is reduced from 35 tokens to 34 tokens as a  
5 token is consumed. Also, the initial packet for the new connection is to be routed  
6 to application endpoint IP2.

7 At (5A), the initial packet is forwarded from classifier 304 to application  
8 endpoint IP2 of host 108(2). Before, during, or after this forwarding, classifier  
9 304 at (5B) plumbs a route for this connection in local DAM table portion 2206.  
10 Classifier 304 may also add an atomic entry 304 for this session into DAM table  
11 2206 for distribution and replication purposes. At (6), future packets for this  
12 connection/session are forwarded from forwarder 302 to application endpoint IP2  
13 of host 108(2) using the local routing table of forwarder 302 as realized by local  
14 DAM table portion 2206 in FIG. 30.

15 FIG. 31 illustrates exemplary high availability mechanisms for network  
16 load balancing infrastructure 106. Specifically, exemplary failure detection 3104,  
17 exemplary failure handling 3106, and exemplary failure recovery 3108 are shown.  
18 These exemplary high availability mechanisms are described with regard to  
19 different network load balancing infrastructure 106 components. The network  
20 load balancing infrastructure 106 components include a forwarder 302, a classifier  
21 304, a request router 306, a session tracker 308, and a health and load handler 314.

22 At 3102(A), forwarder 302 undergoes a local failure. At 3104(A), at least  
23 one load-balancing-aware switch detects the failure. To handle local failure  
24 3102(A), packets are redirected to other forwarder(s) at 3106(A) by the load-  
25 balancing-aware switch. To recover from the failure of forwarder 302, routes that

1 were stored locally at forwarder 302 are rebuilt at 3108(A) at the forwarder(s) to  
2 which packets are redirected using a distributed session tracking manager and a  
3 table thereof such as a DAM and a DAM table thereof. The distributed session  
4 tracking manager may therefore include data redundancies of one or more levels.

5 At 3102(B), classifier 304 undergoes a local failure. At 3104(B), at least  
6 one forwarder detects the failure. To handle local failure 3102(B), packets are  
7 redirected to other classifier(s) at 3106(B) by the forwarder detecting the failure.  
8 To recover from the failure of classifier 304, session information that was stored  
9 locally at classifier 304 are rebuilt at 3108(B) at the classifier(s) to which packets  
10 are redirected using DAM. This session information may be, for example, session  
11 information of a higher level than baseline TCP/IP connections. Also, such  
12 session information may be considered as part of session tracking infrastructure  
13 that is resident on the same device as classifier 304.

14 At 3102(C), request router 306 undergoes a local failure. At 3104(C), at  
15 least one forwarder and/or load-balancing-aware switch detect the failure. To  
16 handle local failure 3102(C), packets are redirected to other request router(s) at  
17 3106(C) by the forwarder and/or load-balancing-aware switch. Individual current  
18 logical requests on which request router 306 is working upon the occurrence of  
19 local failure 3102(C) may be lost unless each such individual logical request is  
20 replicated while the request is being serviced. To recover from the failure of  
21 request router 306, session information and/or routes that were stored locally at  
22 request router 306 are rebuilt at 3108(C) at the request router(s) to which packets  
23 (and thus new logical requests) are redirected. The session information rebuilding  
24 may be effectuated using DAM. Again, such session information may be  
25

1 considered as part of session tracking infrastructure that is resident on the same  
2 device as request router 306.

3 At 3102(D), session tracker 308 undergoes a local failure. At 3104(D), at  
4 least one forwarder and/or classifier detect the failure. For example, if session  
5 tracker 308 is resident on a same device as a classifier, then a forwarder or another  
6 classifier may detect the failure. If session tracker 308 is resident on a separate  
7 device, then a classifier may detect the failure. To handle local failure 3102(D),  
8 data redundancy of one or more levels and distribution across multiple devices are  
9 instituted at 3106(D) for the tracked session information. It should be noted that  
10 the redundancy and distribution are instituted prior to failure 3102(D). To recover  
11 from the failure of session tracker 308, session information from the tables of the  
12 DAM may be redistributed and re-replicated at 3108(D) across at least two devices  
13 (if not already so distributed and sufficiently replicated) in order to handle a  
14 second level of failure.

15 At 3102(E), health and load handler 314 undergoes a local failure. At  
16 3104(E), at least one classifier and/or request router detect the failure. For  
17 example, a component that is receiving health and load information from health  
18 and load handler 314 may detect a failure if health and load handler 314 becomes  
19 non-responsive, especially if health and load handler 314 is resident on a different  
20 device from that of the inquiring component. To handle local failure 3102(E),  
21 cached health and load data redundancy and intrinsic failure handling are  
22 employed at 3106(E) for the health and load information.

23 For example, each health and load handler 314 can include a consolidated  
24 health and load information cache 1208 that duplicates information in health and  
25 load tables 1204 on multiple hosts 108. Also, consumers of the health and load

1 information 1206 of a given health and load handler 314 may be located on a same  
2 device as health and load handler 314 so that failure of health and load handler  
3 314 is intrinsically acceptable. Similarly, the authoritative version of a respective  
4 portion of health and load information 1206 is located on a respective host 108 so  
5 that failure of the host 108 renders the loss of the respective portion of the health  
6 and load information acceptable.

7 To recover from the failure of health and load handler 314, a given network  
8 load balancing component that consumes health and load information may query a  
9 different health and load handler because each such health and load handler  
10 includes a consolidated cache of health and load handler information. Also, when  
11 health and load handler 314 is again accessible, message protocol 1500 may be  
12 used at 3108(E) to rebuild its consolidated cache of health and load information.  
13 Using these exemplary high availability mechanisms, failures of network load  
14 balancing infrastructure 106 components can be detected, handled, and recovered  
15 from in order to mask such failures from clients 102.

#### 16 **Exemplary Connection Migrating**

##### 17 **with Optional Tunneling and/or Application-Level Load Balancing**

18 This section describes how connection manipulation, such as connection  
19 migration, may be utilized in network load balancing. This section primarily  
20 references FIGS. 32-39 and illuminates connection migrating functionality such as  
21 that provided by connection migrator 310 (of FIG. 3). As described above with  
22 reference to FIGS. 3 and 4, each incoming connection at load balancing  
23 infrastructure 106 may be terminated thereat. Afterwards, the connection may be  
24 migrated to a host 108 such that the connection is then terminated at the host 108.  
25 Connection migrator 310 is capable of performing this connection migration and

1 may be located partially at hosts 108 to effectuate the migration. Such connection  
2 migration may be performed in conjunction with application-level load balancing  
3 by a classifier 304 and/or using tunneling via tunneler 312.

4 FIG. 32 illustrates an exemplary approach to application-level network load  
5 balancing with connection migration. Application-level, or layer-7, load balancing  
6 pertains to making load balancing decisions with regard to an application that is to  
7 handle a connection. To perform application-level load balancing, load balancing  
8 infrastructure 106 usually takes into consideration a data portion of a connection.  
9 Unless request routing is employed, a classifier 304 typically takes a peek at the  
10 initial portion of a connection and then migrates the connection, in conjunction  
11 with connection migrator 310, to a selected host 108.

12 For application-level load balancing in a TCP-based environment generally,  
13 classifiers 304 peek at the initial portion of a client's TCP data when deciding  
14 where to forward the client's TCP connection. Thus, application-level logic  
15 examines the client's data and makes load balancing decisions based on that data.  
16 For example, if a connection is an (unencrypted) HTTP connection, a classifier  
17 304 can take a peek at the HTTP header of the first HTTP request in the  
18 connection, and it can make routing decisions based on some portion of the  
19 content of the header (e.g., the URL, a cookie, etc.). Although application-level  
20 load balancing, connection migration, and tunneling are applicable to other  
21 protocols, TCP/IP is used predominantly in the examples herein.

22 As illustrated, load balancing infrastructure 106 (not specifically indicated)  
23 includes a forwarder 302, a classifier 304, a tunneler 312, and a connection  
24 migrator 310 (and possibly e.g. load-balancing-aware router/switches 202(LBA)).  
25 Forwarder 302 corresponds to the virtual IP address and forwards packets to hosts



1 108 in accordance with host selections by classifier 304. Although not specifically  
2 shown in FIG. 32 for clarity, hosts 108 also include connection migrator 310  
3 functionality and tunneler 312 functionality.

4 In a described implementation, forwarder 302, classifier 304, and  
5 connection migrator 310 (at classifier 304 and on hosts 108), along with TCP  
6 protocol software on classifier 304 and hosts 108, cooperate to provide connection  
7 migration. The connection migration illustrated in FIG. 32 is for a connection  
8 from client 102(1) that is initially terminated at classifier 304. After connection  
9 migration, the connection from client 102(1) is terminated at host 108(1). Once  
10 the connection is terminated at host 108(1), packets for the connection may be  
11 tunneled using tunneler 312 (at forwarder 302 and host 108(1)).

12 At (1), client 102(1) sends a SYN packet to forwarder 302 to signal the start  
13 of a new TCP connection. At (2), forwarder 302 forwards this packet to classifier  
14 304. At (3), classifier 304 accepts the TCP connection on behalf of a host 108  
15 (whose identity is not yet known because the actual target host 108( ) has yet to be  
16 selected). In TCP protocol terms, classifier 304 sends a SYN-ACK packet to  
17 client 102(1).

18 At (4), client 102(1) begins sending data. (The initial SYN packet may also  
19 contain data.) The data is processed by classifier 304, which can consult  
20 application-specific logic. The application-specific logic can relate to which host  
21 108 is capable of handling or best handling which types of requests or  
22 connections. Hence, classifier 304 uses the data, as well as application health and  
23 load information from health and load handler 314 and optionally application  
24 session information from session tracker 308, to determine a host 108 that is better  
25

1 or best suited to handle this connection from client 102(1). In this example, host  
2 108(1) is selected.

3 At (5), classifier 304 sends a “binary blob” that represents the state of the  
4 TCP connection to host 108(1). This connection state is aggregated with  
5 cooperation from a TCP stack on classifier 304 by connection migrator 310. The  
6 binary blob contains data from client 102(1) that has been acknowledged by  
7 classifier 304 and TCP parameters such as the TCP/IP 4-tuple, initial sequence  
8 numbers, and so forth.

9 At (6), a connection migrator 310 component on host 108(1) (not explicitly  
10 shown in FIG. 32) “injects” this connection into a TCP stack on host 108(1) using  
11 the state of the TCP connection from the binary blob received from classifier 304.  
12 This connection state injection is performed in cooperation with the TCP stack on  
13 host 108(1), making it appear to applications 316 on host 108(1) that this  
14 connection was originally accepted by host 108(1) itself. Client 102(1) and  
15 applications 316 on host 108(1) are unaware of the connection migration.

16 At (7), classifier 304, in cooperation with the TCP stack on classifier 304,  
17 cleans up the internal state maintained for this connection. This internal state  
18 cleanup at classifier 304 is performed silently such that client 102(1) is not  
19 notified that the connection state is being torn down. Classifier 304 also adds a  
20 route in a local routing table of forwarder 302 that indicates host 108(1) as the  
21 destination for packets of this connection.

22 At (8), subsequent packets for the connection are routed by forwarder 302  
23 to host 108(1) without diversion to or through classifier 304. These packets may  
24 be treated the same by forwarder 302 as those packets for connections that are  
25 classified and routed without using connection migration. These subsequent

1 packets may optionally be tunneled from forwarder 302 to host 108(1) using  
2 tunneler 312. Tunneler 312 is also illustrated (using dashed lines) at connection  
3 migrator 310 at classifier 304 because certain parameter(s) used by tunneler 312  
4 may be determined during a connection migration and/or associated with a  
5 connection being migrated. Exemplary implementations for tunneler 312 are  
6 described further below with particular reference to FIGS. 38 and 39.

7 FIG. 33 is a flow diagram 3300 that illustrates an exemplary method for  
8 migrating a connection from a first device to a second device. Flow diagram 3300  
9 includes seven blocks 3302-3314. Although FIGS. 32 and 34-37 focus primarily  
10 on connection migration in a network load balancing environment, connection  
11 migration as described herein may be effectuated between two devices in general  
12 that each include connection migration functionality, such as that of connection  
13 migrator 310.

14 At block 3302, a connection is accepted at a first device. For example, a  
15 first device may terminate an incoming connection in accordance with one or more  
16 protocols of a protocol stack portion of a network stack. At block 3304, data is  
17 received for the connection at the first device. For example, this data may be  
18 received in an initial packet that requests the connection or in one or more packets  
19 that are received subsequent to an acceptance of the connection.

20 At block 3306, a connection state for the accepted connection is aggregated  
21 from a protocol stack (or more generally from a network stack) at the first device.  
22 For example, a protocol state of the one or more protocols of the protocol stack  
23 may be compiled and aggregated with any received data that has been  
24 acknowledged. At block 3308, the connection state is sent from the first device to  
25

1 a second device. For example, the aggregated information of the connection state  
2 may be sent using a reliable protocol to a second device.

3 At block 3310, the connection state for the connection being migrated is  
4 received from the first device at the second device. At block 3312, the connection  
5 state is injected into a protocol stack (or more generally into the network stack) of  
6 the second device. For example, the connection may be rehydrated using the  
7 protocols of the protocol stack of the second device such that programs above the  
8 protocol stack level are unaware that the connection is a migrated connection.  
9 More specifically, the protocol state may be infused into the protocol stack. The  
10 aggregated data of the connection state is also incorporated at the second device.  
11 At block 3314, the connection is continued at the second device. For example, the  
12 connection may be continued at the second device as if the connection was not  
13 previously terminated elsewhere.

14 FIG. 34 illustrates an exemplary approach to connection migration from the  
15 perspective of an originating device 3400. Connection migration in originating  
16 device 3400 is effectuated, at least partly, by connection migrator 310. In a  
17 described implementation, originating device 3400 is a device that is part of  
18 network load balancing infrastructure 106. For example, originating device 3400  
19 may comprise a classifier 304, possibly along with a forwarder 302, a request  
20 router 306, and so forth.

21 As illustrated, originating device 3400 includes as parts of its network stack  
22 a physical network interface (PNI) 3410, a PNI miniport 3408, a protocol-  
23 hardware interface 3406, a protocol stack 3404, and a socket layer 3402.  
24 Originating device 3400 also includes load balancing functionality 106, such as a  
25 classifier 304 at an application level and connection migrator 310. Specifically,

1 connection migrator 310 includes a migrator intermediate driver 3414 and a  
2 migrator shim 3412. Connection migrator 310 is capable of offloading a  
3 connection from originating device 3400.

4 In a described implementation, physical network interface 3410 may be a  
5 network interface card (NIC) (e.g., an Ethernet NIC), a wireless interface, and so  
6 forth. Although only one physical network interface 3410 is shown, a given  
7 device may actually have multiple such physical network interfaces 3410 (i.e.,  
8 originating device 3400 may be multi-homed). Each physical network interface  
9 3410 typically corresponds to one or more physical network addresses.

10 PNI miniport 3408 is a software module that understands and interfaces  
11 with the specific hardware realization of physical network interface 3410.  
12 Protocol-hardware interface 3406 is a layer that includes one or more respective  
13 interfaces between one or more respective protocols and PNI miniport 3408.

14 Protocol stack 3404 includes one or more respective modules that are each  
15 directed to one or more respective protocols. Examples of such protocols are  
16 described further below with reference to FIGS. 36 and 37. In a transient context,  
17 protocol stack 3404 includes a protocol state 3420 for each connection existing at  
18 originating device 3400. A socket layer 3402 lies between a program such as load  
19 balancing functionality 106 and protocol stack 3404. Socket layer 3402 provides  
20 APIs between load balancing functionality 106 and protocol stack 3404, and it  
21 enables programs to register for connections, among other things.

22 Migrator intermediate driver 3414, or more generally migrator driver 3414,  
23 is located at protocol-hardware interface layer 3406. Migrator shim 3412 is  
24 located transparently between protocol stack 3404 and socket layer 3402.  
25

1       When an initial packet (not shown) requesting a new connection is  
2       presented to originating device 3400, the packet is directed upward from physical  
3       network interface 3410, to PNI miniport 3408, through protocol-hardware  
4       interface layer 3406, and to protocol stack 3404. As the packet traverses the one  
5       or more protocols of protocol stack 3404, protocol state 3420 is created thereat.  
6       Also, as a result of this initial packet or as a consequence of load balancing  
7       functionality 106 accepting the connection to take a peek at the request, data 3416  
8       arrives at originating device 3400.

9       In operation, migrator intermediate driver 3414 diverts a copy of data 3416  
10      to the logic of connection migrator 310. When load balancing functionality 106  
11      issues a migrate connection function call, the migrate function call is passed to a  
12      topmost layer of protocol stack 3404 so that connection state aggregation 3418  
13      may commence. Protocol state 3420 is compiled from the one or more protocols  
14      of protocol stack 3404. In a TCP/IP implementation, protocol state 3420 may  
15      include (i) destination and source TCP ports and IP addresses (e.g., a TCP/IP 4-  
16      tuple), (ii) TCP window state, (iii) initial sequence numbers, (iv) timeout  
17      information, (v) IP fragment ID, (vi) routing information, and (vii) so forth.

18      Connection state aggregation 3418 also aggregates data 3416 that has been  
19      diverted to connection migrator 310 and that has already been acknowledged from  
20      originating device 3400 (e.g., by load balancing functionality 106). This  
21      aggregated connection state 3418 includes protocol state 3420 and data 3416 (and  
22      optionally other connection-related information). Aggregated connection state  
23      3418 is then sent as a binary blob 3422 away from originating device 3400 toward  
24      a targeted device.

1 Binary blob 3422 may be sent from originating device 3400 toward a  
2 targeted device using a reliable protocol. "Reliable" may imply, for example, that  
3 binary blob 3422 is received intact at the targeted device even if individual packets  
4 that constitute binary blob 3422 are lost or corrupted. This binary blob 3422 may  
5 also be bundled with a flow identifier if the connection is to be tunneled  
6 subsequently with tunneler 312. Flow identifiers with tunneling are described  
7 further below with particular reference to FIGS. 38 and 39.

8 FIG. 35 illustrates an exemplary approach to connection migration from the  
9 perspective of a target device 3500. Target device 3500 is similar to originating  
10 device 3400 with respect to the various illustrated layers/modules, including  
11 connection migrator 310. As illustrated however, at least one application 316 at an  
12 application level is interfacing with socket layer 3402. Target device 3500 may  
13 therefore comprise a host 108. Also, connection migrator 310 is capable of  
14 uploading a connection from originating device 3400.

15 In a described implementation, application 316 is the destination of the  
16 connection-initiating packet received at originating device 3400. From originating  
17 device 3400, target device 3500 receives binary blob 3422. Binary blob 3422  
18 includes the connection state associated with the connection being migrated to  
19 target device 3500 and optionally a flow identifier. This connection state includes  
20 protocol state 3420 and acknowledged data 3416 (and possibly other connection-  
21 related information).

22 In operation, when binary blob 3422 reaches protocol-hardware interface  
23 layer 3406, migrator intermediate driver 3414 recognizes it as a blob for  
24 connection migration and diverts it. The connection state is injected at 3502 to  
25

1 create the appearance to application 316 that the connection was originally  
2 terminated at target device 3500.

3 Specifically, protocol state 3420 of injected connection state 3502 is infused  
4 into protocol stack 3404. In a described implementation, protocol state 3420 is  
5 infused first at higher-level protocols and then at lower-level protocols of protocol  
6 stack 3404. After protocol state 3420 is infused into protocol stack 3404, data  
7 3416 can be indicated up to application 316. This data 3416 can be provided to  
8 application 316 as if it were part of a newly and locally terminated connection.

9 After connection state injection 3502 is completed, the connection initiated  
10 by the packet received at originating device 3400 is successfully migrated  
11 therefrom to target device 3500. Subsequent packets for the connection may be  
12 forwarded directly to target device 3500 without passing through originating  
13 device 3400, or at least with only simple routing and no application-level analysis  
14 being applied thereto. Optionally, these packets may be tunneled such that  
15 migrator intermediate driver 3414 effectively operates as a software-based virtual  
16 NIC that is bound to the virtual IP address. In other words, migrator intermediate  
17 driver 3414 (of FIG. 35) may comprise a virtual network adapter that is bound to  
18 the destination address of un-encapsulated packets.

19 FIG. 36 illustrates an exemplary approach to an offloading procedure 3600  
20 for a connection migration. Migration offloading procedure 3600 illustrates  
21 additional exemplary details for a connection migration by an originating device  
22 3400. As illustrated, general protocol stack 3404 includes a TCP stack 3404(T), an  
23 IP stack 3404(I), and an address resolution protocol (ARP) stack 3404(A).  
24 However, other specific protocol stacks 3404( ) may alternatively be employed.  
25



1 By way of example, protocol-hardware interface layer 3406 may be  
2 realized as a network driver interface specification (NDIS)-based layer in a  
3 Microsoft® Windows® operating system (OS) environment. Also, socket layer  
4 3402 may be realized as a Winsock™ layer in a Microsoft® Windows® OS  
5 environment.

6 In a described implementation, migrator intermediate driver 3414 includes  
7 protocol-hardware interfaces 3406 at the junctions to ARP stack 3404(A) and to  
8 PNI miniport 3408. Migrator intermediate driver 3414 serves as an offload target  
9 in migration offloading procedure 3600. The offload target is a protocol-hardware  
10 interface 3406 miniport as illustrated in this example. In a migration uploading  
11 procedure 3700 (as in FIG. 37), migrator intermediate driver 3414 serves as an  
12 upload diverter.

13 More specifically, migrator intermediate driver 3414 is bound to each  
14 physical network interface 3410 through which a TCP connection may be  
15 migrated. Migrator intermediate driver 3414 usually operates as a pass-through  
16 driver by passing packets upwards or downwards in the network stack without  
17 otherwise interacting with the packets. However, migrator intermediate driver  
18 3414 does interact with packets related to connection migration (optionally  
19 including subsequently tunneled packets).

20 Responsibilities of migrator intermediate driver 3414 include: (i) the  
21 acceptance of migrate offload requests; (ii) the aggregation of the protocol state  
22 information that is related to the TCP connection being migrated as compiled from  
23 the specific protocol stacks 3404( ), along with acknowledged data to produce the  
24 connection state information; and (iii) the transmission of the aggregated  
25 connection state to a targeted device 3500 for a migration uploading procedure

1 3700. A reliable wire protocol for such transmission may be shared with that used  
2 by the session tracking components 2002 and 2010 to send and receive session  
3 information messages 2008 (e.g., as described above with reference to FIG. 20).

4 Another responsibility of migrator intermediate driver 3414 (e.g., in a  
5 migration uploading procedure 3700) is to initiate the uploading of migrated  
6 connections that it receives from other devices and to buffer any incoming packets  
7 related to the migrating connection while it is in the process of being uploaded. To  
8 upload the connection, migrator intermediate driver 3414 sends an upload request  
9 to migrator shim 3412. Migrator shim 3412 issues an inject call down into  
10 protocol stack 3404 at TCP stack 3404(A) to instantiate the connection in the  
11 protocol stack 3404 portion of the network stack.

12 Migrator shim 3412 exposes a transport layer client interface to TCP stack  
13 3404(T) and exposes a transport layer provider interface to socket layer 3402.  
14 Migrator shim 3412 has two roles: (i) to initiate connection migration offload  
15 procedure 3600 on an originating device 3400 and subsequently migration upload  
16 procedure 3700 on a targeted device 3500 and (ii) to mediate the classification  
17 process between a host application 316 program, a load-balancing classifier 304  
18 program, and socket layer 3402. Migrator shim 3412 and migrator intermediate  
19 driver 3414 are both further described below with reference to FIGS. 36 and 37.

20 For an exemplary migration offloading procedure 3600, the migration of a  
21 TCP connection is performed after classifier 304 classifies the incoming TCP  
22 connection using one, two, or more packets thereof. Migration offloading  
23 procedure 3600 is described at points <1> through <7>.

24 At <1>, an initialization is performed prior to classification operations.  
25 Protocol stack 3404 makes queries at protocol-hardware interface layer 3406 to

1 determine what offloading capabilities, if any, are available. Migrator  
2 intermediate driver 3414 indicates that connection migration offloading is  
3 available and propagates the query down to PNI miniport 3408. If a TCP chimney  
4 offload ability is provided by a physical network interface 3410, PNI miniport  
5 3408 also so indicates. TCP chimney offload enables some TCP/IP processing to  
6 be offloaded to the hardware of physical network interface 3410 and involves  
7 some compiling of protocol state 3420. Consequently, some compiling and  
8 aggregation logic may be shared between the two offloading mechanisms.

9 At <2>, once a TCP connection has been classified, classifier 304 initiates a  
10 TCP connection migration to a selected host 108. Specifically, a migration  
11 command indicating a targeted device 3500 is issued via socket layer 3402 to  
12 migrator shim 3412.

13 At <3>, migrator shim 3412 initiates TCP connection migration to compile  
14 the TCP protocol state. Specifically, migrator shim 3412 invokes a TCP initiate  
15 migrate offload API (or more generally a migrate connection function call or  
16 migrate connection command). This routine compiles the relevant state for the  
17 specified TCP connection that is used to reinstate the connection on the targeted  
18 device 3500. The compiled protocol state 3420 includes state from the  
19 intermediate stack layers, including TCP stack 3404(T), IP stack 3404(I), and ARP  
20 stack 3404(A).

21 At <4>, once protocol stack 3404 has compiled protocol state 3420 for the  
22 TCP connection being migrated, it invokes an initiate migrate offload API on the  
23 miniport to which it is bound; in this example, that miniport is migrator  
24 intermediate driver 3414. However, in practice, there may be other intermediate  
25 drivers inserted between protocol stack 3404 and migrator intermediate driver

1 3414, such as IP QoS. If so, those IM drivers may participate in the migration, if  
2 relevant, by compiling/aggregating their state to the connection state information  
3 for the connection being migrated. Intermediate drivers continue to propagate the  
4 initiate migrate offload call down the network stack, which eventually results in  
5 execution of a migrate offload handler at migrator intermediate driver 3414. At  
6 this point, migrator intermediate driver 3414 also aggregates any acknowledged  
7 data with the remaining connection state for transfer of the TCP connection to  
8 targeted device 3500.

9 At <5>, after storing/copying connection state information for the TCP  
10 connection being migrated, migrator intermediate driver 3414 notifies the network  
11 stack that the migration is in its final stages by invoking an initiate migrate offload  
12 complete API. This initiate migrate offload complete API follows the reverse path  
13 up the network stack, through the same intermediate drivers (if any), and  
14 eventually to protocol stack 3404. As each layer processes this call, state  
15 information that is associated with the migrated connection may be released. Until  
16 the processing of this call is complete, each layer may send updating notifications  
17 down the network stack to update any part of the connection state that has changed  
18 since the migration was initiated.

19 At <6>, when the initiate migrate offload complete routine reaches TCP  
20 stack 3404(T), TCP silently (i.e., no reset is sent to client 108) closes the  
21 connection, flushing all state associated with the migrated connection, and  
22 propagates the initiate migrate offload complete call to migrator shim 3412. At  
23 this point, the network stack is free of any residual knowledge of the migrated  
24 TCP connection.  
25

1       At <7>, when the initiate migrate offload complete call returns to migrator  
2 intermediate driver 3414 (via the migrator shim 3412 portion of connection  
3 migrator 310), the migration of the TCP connection from originating device 3400  
4 to targeted device 3500 may commence with the transfer of the connection state  
5 thereto. The connection state may be transferred asynchronously and reliably.

6       Once migration is initiated, originating device 3400 is also responsible for  
7 ensuring that subsequent data from client 108 is forwarded to target device 3500.  
8 Consequently, even after the connection is successfully migrated to the target, the  
9 originator retains some amount of state for the connection (e.g., a routing table  
10 entry) in order to properly route subsequent packets to the target. When the  
11 connection is terminated, the target notifies the originator to enable it to purge  
12 whatever residual state remains for the migrated connection.

13       Furthermore, as a consequence of the asynchronous nature of the  
14 connection migration, data packets for the migrating connection that are forwarded  
15 by originating device 3400 (or a forwarder designated thereby if a separate device)  
16 may start arriving at targeted device 3500 before targeted device 3500 receives the  
17 migrated connection state. Migrator intermediate driver 3414 at targeted device  
18 3500 is responsible for buffering those packets until the associated migrated  
19 connection is established on targeted device 3500.

20       FIG. 37 illustrates an exemplary approach to an uploading procedure 3700  
21 for a connection migration. Migration uploading procedure 3700 illustrates  
22 additional exemplary details for a connection migration by targeted device 3500.

23       When a migrated connection arrives at targeted device 3500, it is relayed to  
24 migrator intermediate driver 3414 for processing. After amalgamating and  
25 assimilating the migrated connection state, migrator intermediate driver 3414, in

1 conjunction with migrator shim 3412, injects the migrated connection into the  
2 local network stack in a manner transparent to application 316. For an exemplary  
3 migration uploading procedure 3700, the migration of a TCP connection at points  
4 <1> through <8> is described.

5 At <1>, as described above with reference to migration offloading  
6 procedure 3600, an initialization is performed prior to application hosting  
7 operations. Specifically, protocol stack 3404 makes queries regarding what  
8 offloading capabilities, if any, are available. Migrator intermediate driver 3414  
9 fills in the TCP connection migration support query to indicate that connection  
10 migration uploading is available and also propagates the query down to PNI  
11 miniport 3408 for possible TCP chimney offload capabilities.

12 At <2>, when connection migration data arrives at target device 3500, the  
13 connection migration information (e.g., a bundled binary blob 3422) is delivered  
14 to migrator intermediate driver 3414. Migrator intermediate driver 3414 re-  
15 assembles the connection state, matches it up with any associated data that has  
16 arrived during the migration, and prepares for the upload onto the network stack.  
17 Any data from client 102 that arrives during the process of uploading the migrated  
18 connection is buffered by migrator intermediate driver 3414. Upon successful  
19 completion of the migration, the data will be delivered to application 316.

20 At <3>, to initiate the upload of the migrated connection into the local  
21 network stack, migrator intermediate driver 3414 notifies migrator shim 3412 that  
22 a migrated connection request has arrived. Migrator intermediate driver 3414 also  
23 delivers the connection state (or at least protocol state 3420) to migrator shim  
24 3412.

1       At <4>, migrator shim 3412 initiates the upload of the migrated connection  
2 by invoking a TCP initiate inject routine (or more generally an infuse protocol  
3 state routine) and by providing the migrated protocol state 3420 to TCP stack  
4 3404(T). At <5>, TCP/IP recreates the migrated connection throughout protocol  
5 stack 3404 using the provided protocol state 3420. This protocol state 3420 may  
6 include one or more of transport state (TCP), path state (IP), neighbor and next-  
7 hop state (ARP), and so forth.

8       At <6>, if the migrated connection is successfully reestablished on target  
9 device 3500, TCP initiates a connect event to a client portion of migrator shim  
10 3412 to indicate that a new connection has been established. There are a multitude  
11 of possible reasons for failure, but common reasons may include the lack of a  
12 corresponding listener, routing failure, etc. In these cases where the network stack  
13 is unable to reestablish the migrated connection, no connect event is indicated and  
14 a failure status is specified in the initiate inject complete call. Connection  
15 migrator 310 is responsible for cleaning up the migration and for sending a reset  
16 notification back to client 102 to abandon the connection.

17       At <7>, migrator shim 3412 acts as a provider to propagate the connect  
18 event to socket layer 3402 so as to indicate to the listening application 316 that a  
19 new connection has been established. If the application 316 accepts the  
20 connection, it processes the requests and responds through normal read and write  
21 socket operations; application 316 can be unaware that the connection was  
22 migrated. If the connection is not accepted by the application 316, TCP terminates  
23 the connection but does not send a reset notification back to client 102. Again, a  
24 failure status is specified in the initiate inject complete call, and connection  
25

1 migrator 310 is responsible for cleaning up the migration and for sending a reset  
2 notification back to client 102 to abandon the connection.

3 A special situation arises when application 316 and classifier 304 are co-  
4 located on the same device: migrator shim 3412 may referee between them.  
5 When both classes of programs reside on the same host 108, they may both be  
6 listening to the same IP address(es) and port(s). However, TCP typically has one  
7 listener per unique IP address and port. Consequently, migrator shim 3412 can  
8 obscure a configuration where two programs are listening on the same IP address  
9 and port by multiplexing the two sockets into a single listener at the TCP layer.

10 In such a case, when connect events arrive at the client portion of migrator  
11 shim 3412, migrator shim 3412 as a provider determines on which listening socket  
12 to deliver the connect notification at socket layer 3402. If there is only one socket  
13 listening to the corresponding IP address and port, then that socket receives the  
14 connect event. If there is more than one socket listening, then the recipient  
15 depends on the context in which the connect event is indicated. If the connect  
16 event is a brand new connection for a virtual IP address, then the connect event is  
17 delivered to classifier 304; if the connect event is for a dedicated IP address (non-  
18 load-balanced IP address) or the result of uploading a migrated connection, then  
19 the connect event is delivered to the target application 316.

20 At <8>, once the injection of the migrated connection is complete, TCP  
21 notifies migrator shim 3412 by invoking the provided initiate inject complete  
22 handler. A status code is provided to notify migrator shim 3412 whether or not the  
23 connection was successfully uploaded. If uploading of the migrated connection  
24 fails, connection migrator 310 is responsible for cleaning up the migration and for  
25 notifying client 102 that the connection has been abandoned by sending it a reset.



1 If the migrated connection was successfully injected into the local network stack,  
2 migrator intermediate driver 3414 may begin delivering any buffered data from  
3 client 102 by passing the received packet(s) up through the packet receive path of  
4 protocol-hardware interface 3406.

5 When a migrated connection is terminated (because uploading failed,  
6 because the migrated connection is subsequently closed through normal means,  
7 etc.), target device 3500 notifies originating device 3400. Originating device 3400  
8 uses these notifications to more efficiently and reliably clean out lingering state for  
9 migrated connections, including routing table entries. Therefore, to account for  
10 successfully migrated connections which terminate arbitrarily in the future,  
11 migrator shim 3412 may monitor their activity and notify migrator intermediate  
12 driver 3414 when the sockets therefor are closed.

13 FIG. 38 illustrates an exemplary approach to packet tunneling between a  
14 forwarder 302 and a host 108. Encapsulated packets 3808 may be tunneled from  
15 forwarder 302 to host 108 without incurring overhead for each packet transmitted.  
16 As described further below, the tunneling is effectuated using a flow identifier  
17 3814 and encapsulation mapping tables 3806 and 3810 of tunnelers 312(F) and  
18 312(H), respectively, of forwarder 302 and host 108, respectively. Flow identifier  
19 3814 is inserted into encapsulated packets 3808.

20 As noted above with reference to FIG. 32, packets for a connection that  
21 arrive subsequent to a connection migration may be routed by forwarder 302 to  
22 host 108(1) using tunneling by a tunneler 312. At (8) (of FIG. 32), forwarder 302  
23 forwards such subsequent packets from forwarder 302 having a network address  
24 of "F" to host 108(1) having a network address of "H1". As described above with  
25

1 reference to FIG. 4, forwarder 302 may perform NAT, half-NAT, tunneling, etc. in  
2 order to route the incoming packets to host 108(1).

3 Such incoming packets include a destination IP address of the virtual IP  
4 (“VIP”) address and a source IP address of “C1” for packets arriving from client  
5 102(1). The packets being routed to host 108(1) have a destination IP address of  
6 H1 and a source address of C1 (for half-NAT) or “F” (for full NAT). This re-  
7 writing of the addresses can interfere with some protocols that expect both of  
8 client 102(1) and host 108(1) to have identical views of the source and destination  
9 addresses.

10 Furthermore, at least with respect to full NAT, return paths from host  
11 108(1) to client 102(1) that do not run through forwarder 302 are prohibitive  
12 because host 108(1) does not know the address of client 102(1). Direct paths from  
13 host 108(1) to client 102(1) are desirable in situations in which traffic from host  
14 108(1) to client 102(1) is especially high and/or significantly greater than traffic in  
15 the opposite direction (e.g., when host 108(1) provides streaming media to client  
16 102(1)).

17 Tunneling by tunnelers 312 as described herein can provide for identical  
18 views with respect to the source and destination addresses (and ports) for clients  
19 102 and applications 316 on hosts 108. By way of example and with reference to  
20 FIGS. 34 and 35, tunneler 312 in each of forwarder 302 and host 108 may operate  
21 as part of or in conjunction with a migrator intermediate driver 3414 of a  
22 connection migrator 310.

23 In a described implementation for FIG. 38, connection migrator 310  
24 provides an encapsulation mapping 3812 between a flow identifier 3814 and a  
25 TCP/IP 4-tuple 3804. Connection migrator 310 may be associated with a classifier

1 304, and connection migrator 310 (optionally along with such a classifier 304)  
2 may be located on a same device as forwarder 302. Alternatively, connection  
3 migrator 310 (as well as the classifier 304) may be located on a different device  
4 from forwarder 302. Encapsulation mapping 3812 may alternatively be provided  
5 by or in conjunction with tunneler 312 functionality that is, for example, located at  
6 and/or associated with a classifier 304.

7 By being mapped to a TCP/IP 4-tuple 3804 in encapsulation mapping 3812,  
8 flow identifier 3814 serves to identify a flow of encapsulated packets 3808 for a  
9 particular connection. TCP/IP 4-tuple 3804 includes network addresses (and ports,  
10 etc.) for the source and destination for a particular connection in accordance with a  
11 TCP/IP protocol, or any similar or analogous protocol. Flow identifier 3814 is 32  
12 bits in a described implementation because this allows the flow identifier to be  
13 encoded in the source and destination port fields of the TCP segment header in the  
14 tunneled packet, which enables the tunneled packet to be transmitted without any  
15 tunneling space overhead. At the destination, the TCP/IP 4-tuple can be  
16 determined by looking up the 4-tuple that is linked to the flow identifier as  
17 extracted from the source and destination port fields. However, flow identifiers  
18 3814 of other lengths may alternatively be used, especially for other protocols  
19 such as internet RTP, etc.

20 Each flow identifier 3814 can identify a unique connection from the device  
21 that is originating the tunneling (which is forwarder 302 in this example). Flow  
22 identifiers 3814 may be generated using any appropriate mechanism, such as an  
23 incrementing connection counter. Alternatively, the TCP/IP receiver Initial  
24 Sequence Number (ISN) generated by the connection migrator can serve as flow  
25 identifiers 3814. Furthermore, TCP/IP 4-tuple 3804 is more generally a

1 source/destination pair. Each source value and destination value of an individual  
2 source/destination pair may include a network node identifier (e.g., network  
3 address, port, some combination thereof, etc.) for the source and destination,  
4 respectively, of a given packet propagating on a particular connection.

5 Connection migrator 310 provides encapsulation mapping 3812 to host 108.  
6 Tunneler 312(H) at host 108 stores encapsulation mapping 3812 in encapsulation  
7 mapping table 3810 as encapsulation mapping entry 3810(1). Tunneler 312(H)  
8 can thereafter use flow identifier 3814 to map to and identify the particular  
9 connection corresponding to TCP/IP 4-tuple 3804. Encapsulation mapping 3812  
10 may optionally be provided to host 108 as part of a bundled binary blob 3422 in a  
11 connection migration operation.

12 Forwarder 302 also includes a tunneler 312(F) component with an  
13 encapsulation mapping table 3806. Encapsulation mapping table 3806 stores an  
14 encapsulation mapping entry 3806(1) that links/maps TCP/IP 4-tuple 3804 for a  
15 particular connection to a flow identifier 3814. Tunneler 312(F) also receives the  
16 mapping information for encapsulation mapping entry 3806(1) from connection  
17 migrator 310 (e.g., as an encapsulation mapping 3812).

18 Although only one encapsulation mapping entry 3806(1) and 3810(1) is  
19 shown, each of encapsulation mapping table 3806 and encapsulation mapping  
20 table 3810 may have multiple such entries. These encapsulation mapping tables  
21 3806 and 3810 may be combined with other information, such as tables for session  
22 information of session tracker 308.

23 When a transmitting device (such as forwarder 302) and a receiving device  
24 (such as host 108) of encapsulated packets 3808 only tunnel between each other,  
25 the encapsulation mapping tables thereof likely have the same encapsulation

1 mapping entries. Otherwise, encapsulation mapping table 3806 and encapsulation  
2 mapping table 3810 likely have a different total set of encapsulation mapping  
3 entries 3806( ) and encapsulation mapping entries 3810( ), respectively.

4 In operation, an incoming packet 3802 for a particular connection is  
5 received at forwarder 302. The particular connection is associated with TCP/IP 4-  
6 tuple 3804. Incoming packet 3802 includes TCP/IP 4-tuple 3804 with a source IP  
7 address (of a client 102), a destination IP address (the virtual IP), a source TCP  
8 port (of the client 102), and a destination TCP port.

9 Tunneler 312(F) accepts incoming packet 3802 for tunneling to host 108.  
10 Using TCP/IP 4-tuple 3804, tunneler 312(F) accesses encapsulation mapping table  
11 3806 to locate encapsulation mapping entry 3806(1). Flow identifier 3814 is  
12 extracted from encapsulation mapping entry 3806(1) as being linked/mapped to  
13 TCP/IP 4-tuple 3804.

14 To create encapsulated packet 3808, tunneler 312(F) inserts flow identifier  
15 3814 into the source and destination port portions of the TCP/IP 4-tuple header.  
16 These two TCP portions are 16 bits each, which allows a 32-bit flow identifier  
17 3814 to be inserted. Also, for the source IP address portion of the TCP/IP 4-tuple  
18 header, tunneler 312(F) inserts the IP address "F" of forwarder 302. For the  
19 destination IP address portion of the TCP/IP 4-tuple header, tunneler 312(F) inserts  
20 the IP address "H" of host 108.

21 Forwarder 302 routes/transmits encapsulated packet 3808 to host 108, and  
22 host 108 receives encapsulated packet 3808 from forwarder 302. The tunneler  
23 312(H) component at host 108 detects that encapsulated packet 3808 is a tunneled  
24 packet that is to be de-encapsulated.

1 Flow identifier 3814 is extracted from encapsulated packet 3808 and used  
2 to look up the corresponding TCP/IP 4-tuple 3804 that is linked thereto in  
3 encapsulation mapping entry 3810(1) of encapsulation mapping table 3810.  
4 TCP/IP 4-tuple 3804 is used by tunneler 312(H) to recreate the TCP/IP 4-tuple  
5 3804 header as originally received in incoming packet 3802 at forwarder 302.

6 Specifically, the IP address F of forwarder 302 is replaced with the source  
7 IP address, and the IP address H of host 108 is replaced with the destination IP  
8 address. Furthermore, flow identifier 3814 is replaced by the source TCP port and  
9 the destination TCP port. The de-encapsulated packet is then indicated up the  
10 network stack of host 108 to the targeted application 316.

11 More generally, a portion of a packet header, including a portion of a  
12 source/destination pair, for a given packet that is not necessarily used for  
13 communicating the given packet may be used to carry a flow identifier 3814. By  
14 pre-providing at least part of the source/destination pair at host 108, a flow  
15 identifier 3814 may be employed to tunnel (e.g., encapsulate and/or de-  
16 encapsulate) packets without incurring an encapsulation overhead on each packet.  
17 Furthermore, packets that are full-size with respect to a given protocol may be  
18 tunneled without being fragmented.

19 FIG. 39 is a flow diagram 3900 that illustrates an exemplary method for  
20 packet tunneling between a first device and a second device. For example, the  
21 first device and the second device may correspond to an originating device 3400  
22 and a target device 3500, respectively, of load balancing infrastructure 106 and a  
23 cluster of hosts 108, respectively. Nevertheless, tunneling may be employed in  
24 non-load-balancing implementations.  
25

1 Flow diagram 3900 includes twelve blocks 3902-3924. Although the  
2 actions of flow diagram 3900 may be performed in other environments and with a  
3 variety of software schemes, FIGS. 1-3, 32, 34, 35, and 38 are used in particular to  
4 illustrate certain aspects and examples of the method.

5 At block 3902, a mapping of a flow identifier-to-TCP/IP 4-tuple is sent to a  
6 target device from an originating device. For example, originating device 3400  
7 may send an encapsulation mapping 3812 that links a flow identifier 3814 to a  
8 TCP/IP 4-tuple 3804. At block 3914, the mapping of the flow identifier-to-the  
9 TCP/IP 4-tuple is received at the target device from the originating device. For  
10 example, target device 3500 receives encapsulation mapping 3812 that links flow  
11 identifier 3814 to TCP/IP 4-tuple 3804 from originating device 3400.

12 Alternatively, target device 3500 may receive encapsulation mapping 3812  
13 from another device. As indicated by dashed arrows 3926 and 3928, the actions of  
14 blocks 3904-3912 and blocks 3916-3924 can occur at some time after the actions  
15 of blocks 3902 and 3914, respectively.

16 At block 3904, an incoming packet is received at the originating device  
17 from a client. For example, an incoming packet 3802 having a header with  
18 TCP/IP 4-tuple 3804 may be received at originating device 3400 from a client 102.  
19 At block 3906, a flow identifier is looked up for a connection corresponding to the  
20 client's packet using the TPC/IP 4-tuple of the incoming packet. For example,  
21 flow identifier 3814 may be looked up for the connection with client 102 using  
22 TCP/IP 4-tuple 3804 that is mapped thereto in an encapsulation mapping entry  
23 3806(1) of an encapsulation mapping table 3806.

24 At block 3908, the source IP and destination IP of the incoming packet are  
25 replaced with an originating IP address of the originating device and a target IP

1 address of the target device, respectively. For example, originating device 3400  
2 may replace the IP address portions of the TCP/IP 4-tuple 3804 portion of a header  
3 of incoming packet 3802 with IP addresses of originating device 3400 and target  
4 device 3500.

5 At block 3910, the source port and the destination port of the incoming  
6 packet are replaced with the flow identifier. For example, originating device 3400  
7 may replace source and destination TCP ports of the TCP/IP 4-tuple 3804 portion  
8 of the header of incoming packet 3802 with flow identifier 3814. At block 3912,  
9 the encapsulated packet is sent from the originating device to the target device.  
10 For example, originating device 3400 may send an encapsulated packet 3808 to  
11 target device 3500.

12 At block 3916, the encapsulated packet is received at the target device from  
13 the originating device. For example, target device 3500 may receive the  
14 encapsulated packet 3808 from originating device 3400. At block 3918, the  
15 TCP/IP 4-tuple is looked up for the connection corresponding to the packet  
16 received from the client using the flow identifier. For example, target device 3500  
17 may access an encapsulation mapping table 3810 at an encapsulation mapping  
18 entry 3810(1) that maps flow identifier 3814 to TCP/IP 4-tuple 3804.

19 At block 3920, the originating IP address and the target IP address are  
20 replaced with the source IP address and the destination IP address, respectively,  
21 using the looked-up TCP/IP 4-tuple. For example, target device 3500 may replace  
22 the IP addresses of originating device 3400 and target device 3500 in encapsulated  
23 packet 3808 with the source IP address and the destination IP address from TCP/IP  
24 4-tuple 3804 as attained from encapsulation mapping table 3810.  
25



1       At block 3922, the flow identifier is replaced with the source port and the  
2 destination port of the incoming packet using the looked up TCP/IP 4-tuple. For  
3 example, target device 3500 may replace flow identifier 3814 in encapsulated  
4 packet 3808 with the source TCP port and the destination TCP port from TCP/IP  
5 4-tuple 3804. At block 3924, the client's packet is indicated up to an application at  
6 the target device. For example, a de-encapsulated version of encapsulated packet  
7 3808, or incoming packet 3802, is indicated up to application 316 of target device  
8 3500.

9       The actions, aspects, features, components, etc. of FIGS. 1-39 are  
10 illustrated in diagrams that are divided into multiple blocks. However, the order,  
11 interconnections, layout, etc. in which FIGS. 1-39 are described and/or shown is  
12 not intended to be construed as a limitation, and any number of the blocks can be  
13 combined, rearranged, augmented, omitted, etc. in any manner to implement one  
14 or more systems, methods, devices, procedures, media, APIs, apparatuses,  
15 arrangements, etc. for network load balancing. Furthermore, although the  
16 description herein includes references to specific implementations (and the  
17 exemplary operating environment of FIG. 40), the illustrated and/or described  
18 implementations can be implemented in any suitable hardware, software,  
19 firmware, or combination thereof and using any suitable network organization(s),  
20 transport/communication protocols(s), application programming interface(s)  
21 (APIs), client-server architecture(s), and so forth.

## 22 **Exemplary Operating Environment for Computer or Other Device**

23       FIG. 40 illustrates an exemplary computing (or general device) operating  
24 environment 4000 that is capable of (fully or partially) implementing at least one  
25 system, device, apparatus, component, arrangement, protocol, approach, method,

1 procedure, media, API, some combination thereof, etc. for network load balancing  
2 as described herein. Operating environment 4000 may be utilized in the computer  
3 and network architectures described below or in a stand-alone situation.

4 Exemplary operating environment 4000 is only one example of an  
5 environment and is not intended to suggest any limitation as to the scope of use or  
6 functionality of the applicable device (including computer, network node,  
7 entertainment device, mobile appliance, general electronic device, etc.)  
8 architectures. Neither should operating environment 4000 (or the devices thereof)  
9 be interpreted as having any dependency or requirement relating to any one or to  
10 any combination of components as illustrated in FIG. 40.

11 Additionally, network load balancing may be implemented with numerous  
12 other general purpose or special purpose device (including computing system)  
13 environments or configurations. Examples of well known devices, systems,  
14 environments, and/or configurations that may be suitable for use include, but are  
15 not limited to, personal computers, server computers, thin clients, thick clients,  
16 personal digital assistants (PDAs) or mobile telephones, watches, hand-held or  
17 laptop devices, multiprocessor systems, microprocessor-based systems, set-top  
18 boxes, programmable consumer electronics, video game machines, game consoles,  
19 portable or handheld gaming units, network PCs, minicomputers, mainframe  
20 computers, network nodes, distributed or multi-processing computing  
21 environments that include any of the above systems or devices, some combination  
22 thereof, and so forth.

23 Implementations for network load balancing may be described in the  
24 general context of processor-executable instructions. Generally, processor-  
25 executable instructions include routines, programs, protocols, objects, interfaces,

1 components, data structures, etc. that perform and/or enable particular tasks and/or  
2 implement particular abstract data types. Network load balancing, as described in  
3 certain implementations herein, may also be practiced in distributed processing  
4 environments where tasks are performed by remotely-linked processing devices  
5 that are connected through a communications link and/or network. Especially in a  
6 distributed computing environment, processor-executable instructions may be  
7 located in separate storage media, executed by different processors, and/or  
8 propagated over transmission media.

9 Exemplary operating environment 4000 includes a general-purpose  
10 computing device in the form of a computer 4002, which may comprise any (e.g.,  
11 electronic) device with computing/processing capabilities. The components of  
12 computer 4002 may include, but are not limited to, one or more processors or  
13 processing units 4004, a system memory 4006, and a system bus 4008 that couples  
14 various system components including processor 4004 to system memory 4006.

15 Processors 4004 are not limited by the materials from which they are  
16 formed or the processing mechanisms employed therein. For example, processors  
17 4004 may be comprised of semiconductor(s) and/or transistors (e.g., electronic  
18 integrated circuits (ICs)). In such a context, processor-executable instructions may  
19 be electronically-executable instructions. Alternatively, the mechanisms of or for  
20 processors 4004, and thus of or for computer 4002, may include, but are not  
21 limited to, quantum computing, optical computing, mechanical computing (e.g.,  
22 using nanotechnology), and so forth.

23 System bus 4008 represents one or more of any of many types of wired or  
24 wireless bus structures, including a memory bus or memory controller, a point-to-  
25 point connection, a switching fabric, a peripheral bus, an accelerated graphics port,

1 and a processor or local bus using any of a variety of bus architectures. By way of  
2 example, such architectures may include an Industry Standard Architecture (ISA)  
3 bus, a Micro Channel Architecture (MCA) bus, an Enhanced ISA (EISA) bus, a  
4 Video Electronics Standards Association (VESA) local bus, a Peripheral  
5 Component Interconnects (PCI) bus also known as a Mezzanine bus, some  
6 combination thereof, and so forth.

7 Computer 4002 typically includes a variety of processor-accessible media.  
8 Such media may be any available media that is accessible by computer 4002 or  
9 another (e.g., electronic) device, and it includes both volatile and non-volatile  
10 media, removable and non-removable media, and storage and transmission media.

11 System memory 4006 includes processor-accessible storage media in the  
12 form of volatile memory, such as random access memory (RAM) 4040, and/or  
13 non-volatile memory, such as read only memory (ROM) 4012. A basic  
14 input/output system (BIOS) 4014, containing the basic routines that help to  
15 transfer information between elements within computer 4002, such as during start-  
16 up, is typically stored in ROM 4012. RAM 4010 typically contains data and/or  
17 program modules/instructions that are immediately accessible to and/or being  
18 presently operated on by processing unit 4004.

19 Computer 4002 may also include other removable/non-removable and/or  
20 volatile/non-volatile storage media. By way of example, FIG. 40 illustrates a hard  
21 disk drive or disk drive array 4016 for reading from and writing to a (typically)  
22 non-removable, non-volatile magnetic media (not separately shown); a magnetic  
23 disk drive 4018 for reading from and writing to a (typically) removable, non-  
24 volatile magnetic disk 4020 (e.g., a "floppy disk"); and an optical disk drive 4022  
25 for reading from and/or writing to a (typically) removable, non-volatile optical

1 disk 4024 such as a CD, DVD, or other optical media. Hard disk drive 4016,  
2 magnetic disk drive 4018, and optical disk drive 4022 are each connected to  
3 system bus 4008 by one or more storage media interfaces 4026. Alternatively,  
4 hard disk drive 4016, magnetic disk drive 4018, and optical disk drive 4022 may  
5 be connected to system bus 4008 by one or more other separate or combined  
6 interfaces (not shown).

7 The disk drives and their associated processor-accessible media provide  
8 non-volatile storage of processor-executable instructions, such as data structures,  
9 program modules, and other data for computer 4002. Although exemplary  
10 computer 4002 illustrates a hard disk 4016, a removable magnetic disk 4020, and a  
11 removable optical disk 4024, it is to be appreciated that other types of processor-  
12 accessible media may store instructions that are accessible by a device, such as  
13 magnetic cassettes or other magnetic storage devices, flash memory, compact  
14 disks (CDs), digital versatile disks (DVDs) or other optical storage, RAM, ROM,  
15 electrically-erasable programmable read-only memories (EEPROM), and so forth.  
16 Such media may also include so-called special purpose or hard-wired IC chips. In  
17 other words, any processor-accessible media may be utilized to realize the storage  
18 media of the exemplary operating environment 4000.

19 Any number of program modules (or other units or sets of  
20 instructions/code) may be stored on hard disk 4016, magnetic disk 4020, optical  
21 disk 4024, ROM 4012, and/or RAM 4040, including by way of general example,  
22 an operating system 4028, one or more application programs 4030, other program  
23 modules 4032, and program data 4034.

24 A user may enter commands and/or information into computer 4002 via  
25 input devices such as a keyboard 4036 and a pointing device 4038 (e.g., a

1 “mouse”). Other input devices 4040 (not shown specifically) may include a  
2 microphone, joystick, game pad, satellite dish, serial port, scanner, and/or the like.  
3 These and other input devices are connected to processing unit 4004 via  
4 input/output interfaces 4042 that are coupled to system bus 4008. However, input  
5 devices and/or output devices may instead be connected by other interface and bus  
6 structures, such as a parallel port, a game port, a universal serial bus (USB) port,  
7 an infrared port, an IEEE 1394 (“Firewire”) interface, an IEEE 802.11 wireless  
8 interface, a Bluetooth® wireless interface, and so forth.

9 A monitor/view screen 4044 or other type of display device may also be  
10 connected to system bus 4008 via an interface, such as a video adapter 4046.  
11 Video adapter 4046 (or another component) may be or may include a graphics  
12 card for processing graphics-intensive calculations and for handling demanding  
13 display requirements. Typically, a graphics card includes a graphics processing  
14 unit (GPU), video RAM (VRAM), etc. to facilitate the expeditious display of  
15 graphics and performance of graphics operations. In addition to monitor 4044,  
16 other output peripheral devices may include components such as speakers (not  
17 shown) and a printer 4048, which may be connected to computer 4002 via  
18 input/output interfaces 4042.

19 Computer 4002 may operate in a networked environment using logical  
20 connections to one or more remote computers, such as a remote computing device  
21 4050. By way of example, remote computing device 4050 may be a personal  
22 computer, a portable computer (e.g., laptop computer, tablet computer, PDA,  
23 mobile station, etc.), a palm or pocket-sized computer, a watch, a gaming device, a  
24 server, a router, a network computer, a peer device, another network node, or  
25 another device type as listed above, and so forth. However, remote computing

1 device 4050 is illustrated as a portable computer that may include many or all of  
2 the elements and features described herein with respect to computer 4002.

3 Logical connections between computer 4002 and remote computer 4050 are  
4 depicted as a local area network (LAN) 4052 and a general wide area network  
5 (WAN) 4054. Such networking environments are commonplace in offices,  
6 enterprise-wide computer networks, intranets, the Internet, fixed and mobile  
7 telephone networks, ad-hoc and infrastructure wireless networks, other wireless  
8 networks, gaming networks, some combination thereof, and so forth. Such  
9 networks and communications connections are examples of transmission media.

10 When implemented in a LAN networking environment, computer 4002 is  
11 usually connected to LAN 4052 via a network interface or adapter 4056. When  
12 implemented in a WAN networking environment, computer 4002 typically  
13 includes a modem 4058 or other means for establishing communications over  
14 WAN 4054. Modem 4058, which may be internal or external to computer 4002,  
15 may be connected to system bus 4008 via input/output interfaces 4042 or any  
16 other appropriate mechanism(s). It is to be appreciated that the illustrated network  
17 connections are exemplary and that other means of establishing communication  
18 link(s) between computers 4002 and 4050 may be employed.

19 Furthermore, other hardware that is specifically designed for servers may  
20 be employed. For example, SSL acceleration cards can be used to offload SSL  
21 computations. Additionally, especially in a network load balancing operating  
22 environment, TCP offload hardware and/or packet classifiers on network  
23 interfaces or adapters 4056 (e.g., on network interface cards) may be installed and  
24 used at server devices.

1 In a networked environment, such as that illustrated with operating  
2 environment 4000, program modules or other instructions that are depicted  
3 relative to computer 4002, or portions thereof, may be fully or partially stored in a  
4 remote media storage device. By way of example, remote application programs  
5 4060 reside on a memory component of remote computer 4050 but may be usable  
6 or otherwise accessible via computer 4002. Also, for purposes of illustration,  
7 application programs 4030 and other processor-executable instructions such as  
8 operating system 4028 are illustrated herein as discrete blocks, but it is recognized  
9 that such programs, components, and other instructions reside at various times in  
10 different storage components of computing device 4002 (and/or remote computing  
11 device 4050) and are executed by processor(s) 4004 of computer 4002 (and/or  
12 those of remote computing device 4050).

13 Although systems, media, devices, methods, procedures, apparatuses,  
14 techniques, schemes, approaches, procedures, arrangements, and other  
15 implementations have been described in language specific to structural, logical,  
16 algorithmic, and functional features and/or diagrams, it is to be understood that the  
17 invention defined in the appended claims is not necessarily limited to the specific  
18 features or diagrams described. Rather, the specific features and diagrams are  
19 disclosed as exemplary forms of implementing the claimed invention.  
20  
21  
22  
23  
24  
25